

# Statistisches Modellieren – Zusammenfassung

Einführung – «Daten verlangen keine Gerechtigkeit!» – S. 9, Skript STMO Prof. Dr. Andreas Ruckstuhl  
Statistische Modellierung heisst **mittels Daten ein Modell für die Realität erstellen**. Das Modell soll die **Realität beschreiben, vorhersagen** oder **kausal analysieren**. Mit einem statistischen Modell lassen sich **Genauigkeitsangaben** machen.

Notationen – dat <- read.table("C:/Users/denis/.../Daten.dat", header = T)  
• ( ) = Funktionsargumente • { } = Mengen • ( ) = Priorität der Rechenoptionen • [ ] = Vektoren/Matrizen

Regressionsanalyse – «All models are wrong, but some are useful» – George E. P. Box  
• Ist eine statistische Methode, um die **Beziehung** zwischen (Mess-) **Grössen zu untersuchen** und zu **modellieren**.  
• Modellvorstellung aufbauen: Modell mit unbekanntem Parametern.

Beispiel Sprengungen beim Bau eines Strassentunnel  
Erschütterung der Häuser darf einen bestimmten Wert nicht überschreiten. Erschütterung hängt ab von Ladung, Distanz vom Sprengort, Untergrundmaterial, Ort der Sprengung, etc. Die genaue Funktion zur Berechnung der Grösse der Erschütterung ist unbekannt. Schätzung der Erschütterung aus Daten.

- **Zielgrösse y (response variable)** → die Erschütterung
- Hängt über Funktion **h** von den **erklärenden Variablen**  $x^{(1)}, \dots, x^{(m)}$  (**predictor variables**) ab → Ladung, Distanz, etc.
- Im Idealfall sollte für jede **Beobachtung/Versuch i** gelten  $y_i = h(x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(m)})$ . Diese Formel existiert leider nicht!
- Formel gilt ungefähr, (Abweichungen zufällig) deshalb:  $Y_i = h(x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(m)}) + E_i$  und nennt die  $E_i$  **Zufallsfehler**.
- Vorstellungen über Grösse solcher Abweichungen werden mit einer Wahrscheinlichkeits-Verteilung formuliert.
- Funktion **h( )** beschreibt oft eine Gerade; mit lin. Regression lassen sich auch nicht lineare Kurven an Daten anpassen!

Einfache lineare Regression  
• Besteht nur aus einer erklärenden Variable  $x_i^{(1)}$ . Gerade ist einfachste Funktion, die eine Abhängigkeit ausdrücken kann:  
 $y_i = \alpha + \beta x_i$ , wobei  $\alpha =$  **Achsenabschnitt** (Schnittpunkt mit der y-Achse) und  $\beta =$  **Steigung**  
• modellieren die Abweichungen von der Geraden mit einer Zufallsvariablen:  $Y_i = \alpha + \beta x_i + E_i$   
• Wir nehmen an, dass die Abweichungen  $E_i$   
◦ alle eine bestimmte gleiche Verteilung haben & ◦ stochastisch unabhängig (insbesondere unkorreliert) sind.  
• I.d.R. für Abweichung  $E_i$  **Annahme Normalverteilung** mit  $E(X) = 0$  und  $Var(X) = \sigma^2 \rightarrow E_i \sim N(0, \sigma^2)$   
• Modell erst dann konkret, wenn die Parameter  $\alpha, \beta$  und  $\sigma$  festgelegt sind.  
• Daraus folgt, dass auch die Zielvariable  $Y_i$  eine Zufallsvariable ist, weshalb ein grosser Buchstabe verwendet wird.  
• Zielvariable  $Y_i$  ist Normalverteilt mit konstanter Varianz  $\sigma^2$ , jedoch ist Erwartungswert für jede Beobachtung anders (d.h.  $= \alpha + \beta x_i$ ). Das beobachtete  $y_i$  entsteht dann aus einer Realisation  $e_i$  der Zufallsvariable  $E_i$ , also  $y_i = \alpha + \beta x_i + e_i$

Schätzung der Parameter  
• Die Funktionen, die den Daten die best-passenden Werte zuordnen, heissen Schätzfunktionen oder Schätzungen.  
Herleitung Schätzung nach Prinzip **Kleinsten Quadrate (KQ)**: Parameter  $\alpha, \beta$  so bestimmt, dass **Summe der quadrierten Abweichungen**:  $\sum_{i=1}^n r_i^2$  mit  $r_i = y_i - (\alpha + \beta x_i)$  **minimal** wird.  $\hat{\beta} = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$   
• Die Schätzfunktionen lauten dann  $\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$  und  $\hat{\beta} = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$   
`Spr.lm <- lm(Ersch ~ IDist, data = SprengS1)` Übergabe `data.frame` an Bf lm()  
`coef(Spr.lm)` (Zielgrösse (Y) ~ erkl. Variable (X))  
(Intercept) IDist  
8.9791 -1.9235  
Geschätzte Grade lautet:  $lErsch = 8.9791 + (-1.9235) \cdot IDist$   
`plot(Ersch ~ IDist, data = SprengS1); abline(Spr.lm)` #Linie einzeichnen  
• Intercept = Achsenabschnitt, IDist = Steigung der Geraden (IDist wird von Variablenname des data.frame übernommen).  
• **Weiterer Parameter** im Modell, die Varianz  $\sigma^2 = var(E_i)$  der **zufälligen Abweichungen**. Die zufälligen Fehler  $E_i$  können **weder** direkt beobachtet noch aus  $E_i = Y_i - (\alpha + \beta x_i)$  hergeleitet werden, da  $\alpha$  und  $\beta$  unbekannt sind. Bekannt sind «Näherungswerte» für  $E_i$  (**Residuen**)  $R_i := Y_i - (\hat{\alpha} + \hat{\beta} x_i)$ , erwartungstreue Schätzung für  $\sigma^2$  ist:  $\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n R_i^2$   
• Die Kleinsten-Quadrate-Schätzungen  $\hat{\alpha}$  &  $\hat{\beta}$  sind **erwartungstreu** und **normalverteilt** mit Varianzen  $\sigma^{(\beta)^2} = \frac{\sigma^2}{SS_X}$   
 $\sigma^{(\alpha)^2} = \sigma^2 (\frac{1}{n} + \frac{\bar{x}^2}{SS_X})$  wobei **Quadratsumme**  $SS_X = \sum_{i=1}^n (x_i - \bar{x})^2$ . Sie ist die **effizientesten** Schätzungen, **sofern** die **Zufallsfehler unabhängig** sind und alle die gleiche **Normalverteilung**  $N(0, \sigma^2)$  haben. Treffen diese Voraussetzungen **nicht** zu, kann Methode unzuverlässig werden, womit andere Schätzverfahren besser geeignet sind.

Statistisches Testverfahren (Vorgehen gemäss Modul GSTAT)  
• Sind Daten mit Modell mit (teilweise) vorgegebenen Parametern verträglich?  
• Unter  $H_0: T$  hat eine t-Verteilung mit  $n - 2$  Freiheitsgraden und der Teststatistik:  $T = \frac{\hat{\beta} - \beta_0}{se(\hat{\beta})}$ , mit  $se(\hat{\beta}) = \sqrt{\frac{\hat{\sigma}^2}{SS_X}}$

**Hypothesentest Beispiel Kirschbäume**:  $lV = \log(Volumen) = y$ ,  $lD = \log(Durchmesser) = x_1$ ,  $lH = \log(Höhe) = x_2$   
Modell:  $lV = \beta_0 + \beta_1 lD + \beta_2 lH + E \rightarrow$  Modellidee mit Kegelvolumen Grundlage:  $V = \pi \cdot R^2 \cdot \frac{h}{3} = \pi \cdot (\frac{D}{2})^2 \cdot \frac{H}{3} = \frac{\pi}{(2^2 \cdot 3)} \cdot D^2 \cdot H$

Nun Modellidee in Modellform von oben bringen:  $\log(V) = \log(\frac{\pi}{(2^2 \cdot 3)} \cdot D^2 \cdot H) + 2 \cdot \log(D) + \log(H) = -1.340177 + 2 \cdot \log(D) + \log(H)$   
Testen, ob Modellidee  $\beta_1 = 2$  und  $\beta_2 = 1$  sein kann  $\rightarrow$  Nullhypothese:  $lV = \beta_0 + 2lD + 1lH + E$ ; schätzen  $\beta_0$   
`Chr.lm <- lm(lV ~ lD + lH, data = cherry)`; `Chr.lm0 <- lm(lV ~ offset(2*lD + lH), data = cherry)`; `anova(Chr.lm0, Chr.lm)`  
Vergleich normales Modell mit Modellidee. P-Wert signifikant? Wenn ja, dann haben wir statistische Evidenz gegen Nullhypothese von oben. Wenn nicht signifikant, dann keine Evidenz gegen Nullhypothese.

Vertrauensintervall `summary(Spr.lm)` →  
• Vertrauensintervall umfasst alle **Parameterwerte  $\alpha, \beta$** , die auf Grund bestimmten statistischen Tests nicht abgelehnt werden.  
• Damit Null-Hypothese auf Niveau 5% nicht abgelehnt werden muss, muss Absolutbetrag Testgrösse kleiner als das entsprechende Quantil sein:  $|T| \leq q_{0.975}^{n-2}, \hat{\beta} \pm q \cdot se(\hat{\beta}), se(\hat{\beta}) = \frac{\hat{\sigma}}{\sqrt{SS_X}}$   
• **confint(Spr.lm, level = 0.95)** «Hand»:  $qt(1 - \frac{\alpha}{2}, df = n - 2)$   
 $\alpha = 5\%$  2.5% 97.5% `coef(fit)[2] + c(-1,1)*qt(0.975, fit$df.residual)*summary(fit)$coefficients[2,2]` `Kof=  $\beta_j + 1$`   
(Intercept) 3.204 4.594  $-1.92 + c(-1,1) \cdot 2.200 \cdot 0.1783 = KI$  (Daten v. summary output oben)  
IDist -2.315 -1.531  $KI$  von IDist [-2.315, -1.531], womit Steigung -2 gut mit Daten verträglich.

**Ziel**: Aufgrund Daten etwas über Werte der Parameter des Modells zu sagen, die plausibel erscheinen. Drei Fragen stellen:  
1. **Welcher Wert** ist für den (resp. jeden) Parameter am **plausibelsten**? Antwort wird durch **Schätzung** gegeben  
2. Ist ein **bestimmter Wert plausibel**? Die **Entscheidung** trifft man mit einem **Test**.  
3. Welche **Werte** sind **insgesamt plausibel**? Als Antwort erhält man ein **Vertrauens-/Konfidenzintervall**.

Erwartungswert  $E(Y(x_0))$  für Wert  $x_0$  (**Punktschätzung/Prognose**): `predict(Spr.lm, newdata = data.frame(Dist = 1))`  
**Korrektur Prognosewerte von log. Var**: `(h.p <- predict(Cher.lm, newdata = data.frame(ID = log(5.3), lH = log(27))))`  
`exp(h.p)`; #  $\rightarrow$  ist median anstelle  $E(X)$  versus; `exp(h.p + summary(Cher.lm)$sigma^2/2)`; # ist nun  $E(X)$   
Korrektur bei log Transformation, gilt für Erwartungswert und Vertrauensintervall, aber nicht für Prognoseintervall (dieses ist eher zu lang als zu kurz). siehe auch Abschnitt «Auswirkungen der Transformation von Zielvariablen» in dieser ZF.

Vertrauensintervall VI für  $\eta_0$   
• Oft liegt Interesse an Funktionswert  $h(x_0)$  an bestimmten Stelle  $x_0$ . Bildung Vertrauensintervall.  
•  $h(x_0) = \alpha + \beta x_0 = \eta_0$  (eta). Finden aller plausiblen Werte via eines Tests  $\eta_0 = \hat{\eta}_0$   $se(\hat{\eta}_0) = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_X}}$   
• Schätzung von  $\eta_0$ :  $\hat{\eta}_0 = \hat{\alpha} + \hat{\beta} x_0$ , t-Verteilung mit  $n-2$  Freiheitsgraden Testgrösse  $T = \frac{\hat{\eta}_0 - \eta_0}{se(\hat{\eta}_0)}$   
• Vertrauensintervall für  $\eta_0 = h(x_0)$ :  $(\hat{\alpha} + \hat{\beta} x_0) \pm c_w \cdot se(\hat{\eta}_0)$ , wobei  $c_w$  das Quantil der t-Verteilung mit  $n-2$  Freiheitsgraden  
• **predict(Spr.lm, newdata = data.frame(Dist = c(log(50), log(100))), interval = "confidence", level = 0.95, parm = 2)**  
• Kann Intervalle für mehrere Schätzer (parm = weglassen) gleichzeitig erstellt werden. Variablenname in Data.frame (hier **IDist**) muss mit Name in Datensatz respektive Befehl `Spr.lm <- lm(Ersch ~ IDist, data = SprengS1)` übereinstimmen.

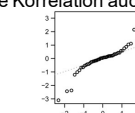
Prognoseintervall  $E_i: \hat{\alpha} + \hat{\beta} x_0 \pm c_w \cdot \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{SS_X}} = \hat{\eta}_0 \pm c_w \cdot \sqrt{\hat{\sigma}^2 + (se^2(\hat{\eta}_0))}$   
• In welchem **Bereich** kommt eine **zukünftige Beobachtung** zu liegen. **Aussage über Zufallsvariable!** Das **Prognoseintervall berücksichtigt Variabilität** von  $E_i \uparrow$   
• **Beispiel Sprengungen**: Wie gross wird die Erschütterung sein, wenn die Distanz zur Sprengstelle 50m/100m beträgt?  
`(h <- predict(Spr.lm, newdata = data.frame(IDist = c(log(50), log(100))), interval = "prediction", level = 0.95))`  
anschliessend **Rücktransformation**, da Werte logarithmiert  
fit lwr upr  
1 1.45428 0.8496141 2.058947 `exp(h)`

• Tabelle für Modell 95%-Prognose-Intervalle für alle Beobachtungen: `cat.p <- predict(cat.lm, interval = "prediction")`  
• `cbind(unten = cat.p[, "lwr"], fit = cat.p[, "fit"], oben = cat.p[, "upr"], laenge = cat.p[, "upr"] - cat.p[, "lwr"])`  
**Plot**: `Forbes.lm <- lm(y ~ x, Forbes[-12,]); x0 <- data.frame(x = seq(min(Forbes$x), max(Forbes$x), length = 50))`  
`Forbes.cia <- predict(Forbes.lm, newdata = x0, interval = "confidence", level = 0.99); plot(Forbes$x, Forbes$y)`  
`lines(x0$x, Forbes.cia[, "upr"], col = 2); lines(x0$x, Forbes.cia[, "lwr"], col = 2); abline(Forbes.lm, col = "blue", lty=1)`  
`Forbes.pia <- predict(Forbes.lm, newdata = x0, interval = "prediction", level = 0.99) #Darstellung in Plot: Erzw. E(X)`  
`lines(x0$x, Forbes.pia[, "upr"], col = 7); lines(x0$x, Forbes.pia[, "lwr"], col = 7) #Vertrauens- und Prognoseband`

Prüfen der Modelleignung  
• Beim **Regressionsmodell beruhen** die eingeführten Schätz- und Testmethoden auf **Modellannahmen**: Die Fehler  $E_i$  sind **unabhängig** und **normalverteilt** mit **konstanter Varianz**,  $E_i$  unabhängig  $\sim N(0, \sigma^2)$ . Annahmen:  
1. Der Erwartungswert der  $E_i$  ist  $E(E_i) = 0$ , 3. sie sind **normalverteilt**.  
2. die  $E_i$  haben alle die gleiche **Varianz**  $var(E_i) = \sigma^2$ , 4. sie sind **unabhängig**.  
• Für die Regressionsfunktion wurde Linearität vorausgesetzt. Wenn Annahme verletzt, gilt für Fehler Annahme (1) nicht.  
• Damit ist klar, dass alle Verletzungen der gemachten Annahmen in den Fehlern  $E_i$  sichtbar sein müssen.  
• Gerade aus Regressionsanpassung in Diagramm legen. Eine Möglichkeit, diese Güte der Anpassung zu quantifizieren, ist das Bestimmtheitsmass  $R^2$ , das im R-Output mit «**Multiple R-squared**» bezeichnet ist.

**Bestimmtheitsmass  $R^2$** : misst Güte der Anpassung mit Anteil der durch die Regression erklärten Streuung der Y-Werte:  
 $R^2 = \frac{\sum_{i=1}^n \hat{y}_i^2}{\sum_{i=1}^n y_i^2} = \frac{SS_{Fit}}{SS_Y}$ ,  $R^2$  liegen zwischen 0 und 1; je grösser desto besser.

• In Technik und Naturwissenschaften: Werte grösser 0.9 durchaus üblich. In Geistes- und Sozialwissenschaften: um 0.6  
• Falls kein Achsenabschnitt  $\beta_0$  im Modell vorhanden, sollte diese Definition des Bestimmtheitsmasses nicht verwenden.  
• **Interpretation**: das Bestimmtheitsmass ist identisch zur  **quadrierten Korrelation** zwischen der Zielvariablen  $Y_i$  und den angepassten Werten (fitted values):  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ . Demzufolge misst das Bestimmtheitsmass wie die Korrelation auch die Stärke des linearen Zusammenhangs zwischen der Zielvariablen und der Anpassung.  
• **Achtung!** Das Bestimmtheitsmass ist kein Mass für die **Eignung des Regressionsmodells!**

**Langschwänzigkeit** →   
19.09.2023 4. Semester Seite 2 von 12

1.) **Tukey-Anscombe-Diagramm (T-A-Diagramm/Plot):**  $\mathbb{E}(E_i) = 0 \rightarrow$  Erwartungswert = 0

- Tragen Residuen  $R_i$  gegen die angepassten Werte  $\hat{y}_i (= \hat{\alpha} + \hat{\beta}x_i)$  auf.
- $\mathbb{E}(E_i) = 0 \rightarrow$  **Residuen** sollten im Tukey-Anscombe-Diagramm in allen Abschnitten um **horizontale Null-Linie streuen**.
- **Idee:** Struktur mit dem **Glätter** «loess» mit lokalen Gerade (anstelle Mittelwert) durch **Bootstrap-Simulationen** unter der Gültigkeit der Modellannahme erzeugen. **Vergleich simulierten Kurven mit ursprünglichen**, ob «extremer»?
- **Vorgehen Bootstrap-Simulationen:** Erzeuge zufällig neue Beobachtungen  $y_i^*$ , die dem Modell entsprechen, d.h. erzeuge  $n \sim \mathcal{N}(0, 1)$ -verteilte Zufallszahlen  $e_i^*$  und bilde daraus neue Zielvariablenwerte  $y_i^* = \hat{y}_i + \hat{\sigma} \cdot e_i^*$ .
- Danach Regressionsrechnung mit erklärenden Variablen durchführen und den neu erzeugten Werten  $y_i^*$  durch, berechnet den Glätter für das T-A-Diagramm und ihn ins Diagramm einzeichnen. Wiederhole dies beiden Schritte  $n = 19$  Mal.
- **Überprüfung in Plot, ob Glättung in Bandbreite der Simulationen liegt und ob Kurvenform untypisch (Banane)!**

2.) **Scale-Location Plot / Streuungs-Diagramm:**  $var(E_i) = \sigma^2 = konst \rightarrow$  konstante gleiche Varianz

- Erfassen Struktur mit einem «gleitenden Streumass». Eine einfache Möglichkeit besteht darin, die zuvor benutzte Methode «loess» auf die **Absolutwerte**  $|R_i|$  der **Residuen** oder besser auf  $\sqrt{|R_i|}$  anzuwenden  $\rightarrow$  **scale-location plot**
- $\sqrt{|R_i|}$  Wurzel verwenden, weil  $\sigma$  auf  $|R_i|$  zu schief verteilt ist. Der Glätter schätzt  $\sigma$  nicht direkt.

3.) **Histogramm/Normal QQ-Plot:**  $E_i \sim \mathcal{N}(0, \hat{\sigma})$  Fehler sind normalverteilt

- Besser als Histogramme ist **QQ-Plot**  $\rightarrow$  Quantile der empirischen Verteilung der Residuen werden mit Quantilen der Normalverteilung verglichen. Falls Daten normalverteilt sind, **streuen Punkte** um eine **Gerade**.
- **Achtung:** Ein normal QQ-plot für die  $Y_i$  ist sinnlos, da die  $Y_i$  ja verschiedene Erwartungswerte haben, weshalb geschätzte  $E_i$  genommen werden, sprich die Residuen. **Plot: streuen (nicht) innerhalb stochastischen Fluktuation?**

`source("../RFN_Plot-ImSim.R"); par(mfrow = c(2, 3)); plot(fit, which = 1:3); plot.lmSim(fit, SEED = 1798, rob = T)`  
**Aussage zu Plots:** Der Glätter liegt ((teilweise) nicht) innerhalb der stochastischen Fluktuation der Simulation / die Punkte streuen gut um eine Gerade und liegen (nicht) innerhalb der stochastischen Fluktuation des QQ-Normalplot. Damit liegt (k)eine statistische Evidenz gegen die Annahme Erwartungswert  $\mathbb{E}(E_i) = 0$  / konstante Varianz  $var(E_i) = \sigma^2$  / Normalverteilung  $E_i \sim \mathcal{N}$  vor. **Frage, ob Voraussetzung der Unabhängigkeit erfüllt ist, ist schwieriger zu beantworten**

**Prüfen der Modelleignung:** wollen mit drei Darstellungen sicherstellen, dass Daten keine für Theorie gefährlichen Abweichungen zu Voraussetzungen gibt. Da Residuen Realisierungen von Zufallsgrößen, Voraussetzungen sind «exakt erfüllt».

**Skalierte Residuen:**

- Wollen die Verteilung der Zufallsfehler  $E_i$  überprüfen, haben aber die Residuen  $R_i$  benützt – und das ist nicht dasselbe!
- Falls Fehler  $E_i$  normalverteilt  $\rightarrow$  Residuen ebenfalls. Aber sie haben **nicht die gleiche Varianz!**
- Damit Residuen gleiche Verteilung haben, muss man sie skalieren:  $\tilde{R}_i = R_i / \sqrt{1 - (\frac{x_i - \bar{x}}{s_{xy}})^2} \sim \mathcal{N}(0, \sigma^2)$
- **Vorgehen:** Zur Überprüfung der Voraussetzungen verwenden der standardisierten Residuen, - ausser im Tukey-Anscombe-Diagramm, weil so Abweichungen von einem konstanten Erwartungswert besser erkannt werden können.

Behandlung von Unzulänglichkeiten: J.W. Tukey nannte sie **First Aid Transformations**

- Annahme «Streuung der Zufallsfehler ist konstant» verletzt. **Häufig:** Abhängigkeit der Streuung von  $\hat{y}$ .
- Die **Streuung nimmt** im Streuungs-Diagramm zu und der **Normal-Plot** zeigt eine **rechts schiefe Verteilung** an.
- Folgende Transformationen können helfen:
  - o **Logarithmus-Transformation für Konzentrationen/Beträge,** o **Wurzeltransformation für Zähdaten**
  - o **Arcus-Sinus-Wurzel-Transformation**  $\hat{y} = \arcsin(\sqrt{\hat{y}})$  oder
  - o **Logit-Transformation für Anteile (Prozentzahlen/100):**  $\hat{y} = \log(\frac{y+0.005}{1.01-y})$  **First-Aid/Log Faktor 10**
- Immer anwenden, ausser triftige Gründe sprechen dagegen. **Nicht anwenden bei Zeitvariablen / zu wenig Streuung**

**Auswirkungen der Transformation von Zielvariablen:** (bei Log Problem für Erwartungswert und Vertrauensintervall)

- Transformation der Zielvariablen ändert Form der Verteilung der Fehler. Beispiel Logarithmus-Transformation:
- Sowohl erklärende Variable als auch Zielgrösse seien Konzentrationen,  $\tilde{Y} = \log(Y)$  und  $\tilde{x} = \log(x)$ .
- Aus  $\tilde{Y}_i = \alpha + \beta \tilde{x}_i + E_i$  wird  $\log(Y_i) = \alpha + \beta \cdot \log(x_i) + E_i$  oder  $Y_i = e^{\alpha} \cdot x_i^{\beta} \cdot e^{E_i}$  **Korrektur**  $\exp(\frac{\hat{\sigma}^2}{2})$ , für  $\hat{y}_0 \rightarrow \exp(\hat{y}_0 + \frac{\hat{\sigma}^2}{2})$
- D.h. haben **Potenzgesetz** für die ursprünglichen Grössen und **Fehler ist proportional** (und **nicht additiv**).
- Falls  $\beta = 1$  ist die Zielvariable proportional zu  $x$  bis auf einen **multiplikativen zufälligen Fehler**.  $\rightarrow$  **Log-Norm-verteilt**
- Bsp:  $\log(B_i) = \beta_0 + \beta_1 \cdot \log(N_i) + \beta_2 \cdot \log(C_i) + \beta_3 \cdot p_i + E_i \Leftrightarrow B_i = e^{\beta_0} \cdot N_i^{\beta_1} \cdot C_i^{\beta_2} \cdot e^{\beta_3 p_i} \cdot e^{E_i}$   $E_i \sim \mathcal{N}(0, \sigma^2)$  unabhängig

**Annahme «Erwartungswert ist konstant null» verletzt:**

- **Systematische Abweichung** im Erwartungswert kann oft durch **Transformation der erklärenden Variablen  $x$  oder durch Hinzufügen eines zusätzlichen Terms  $x^2$  (quadratische Regression) zum Verschwunden** gebracht werden.

**Annahme «Fehler sind normalverteilt» durch Ausreiser verletzt:** *weglassen von Beobachtungen mit Bf subset ↓ ↓*

- Richtigkeit Daten überprüfen! Falls Daten i.O., Transformation der Ziel- und/oder der erklärenden Variablen abklären.

**Annahme «Fehler sind normalverteilt» durch Langschwänzigkeit verletzt:** `lm(y ~ x, Forbes, subset = c(-11, -12))`

- Extremste Beobachtungen weglassen, bis Langschwänzigkeit verschwindet  $\rightarrow$  Resultate mit Vorsicht geniessen, weil zu optimistisch! Kleinste-Quadrate-Methode ist bei Langschwänzigkeit nicht optimal, nur **robuste Methoden geeignet**.

Unabhängigkeit der zufälligen Fehler: Korrelation in der Interpretation

- Wenn Beobachtungen eine zeitliche Reihenfolge einhalten, dann könnten **Autokorrelationen** vorhanden sein.
  - o Also Residuen  $R_i$  in dieser Reihenfolge auftragen – Keine Strukturen sichtbar sein! `plot(resid(fit), type = "h")`
  - o  $R_i$  gegen  $R_{i-1}$  in Streudiagramm auftragen – Punkte sollten frei streuen und **keine Korrelation** zeigen.
  - o `fit <- lm(y ~ x, data = dat); x.n <- length(resid(fit)); scatter.smooth(resid(fit)[1:(x.n - 1)], resid(fit)[2:x.n])`
- **Wenn (Auto-) Korrelationen** vorliegen, dann sind **P-Werte** der üblichen Tests häufig **falsch** und **Vertrauensintervalle zu kurz**. Methoden, die Korrelationen berücksichtigen: **verallgemeinerte Kleinste Quadrate** (R: `gls` aus nlme).

**Multiple lineare Regression**

- **Zusammenhang** zwischen **einer Zielgrösse** und **mehreren erklärenden Variablen**  $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ .  $\rightarrow$  **pairs(dat)**
- Erweiterung des einfachen Regressionsmodell:  $Y_i = \beta_0 + \beta_1 x_i^{(1)} + \beta_2 x_i^{(2)} + \dots + \beta_m x_i^{(m)} + E_i$  mit **unab.  $E_i \sim \mathcal{N}(0, \sigma^2)$**
- Die **(unbekannten) Parameter** sind die **Koeffizienten  $\beta_0, \beta_1, \dots, \beta_m$**  **zu den erklärenden Variablen** und die **Varianz  $\sigma^2$**  der **zufälligen Abweichungen  $E_i$** . Die  $\beta$  sind die **Steigungen in Richtung der x-Achsen**. Achsenabschnitt:  $\beta_0 = \alpha$
- Schätzung Koeffizienten  $\beta_j$  erfolgt via **KQ** und sind **normalverteilt**  $\rightarrow$  Verteilung von Teststatistiken  $\rightarrow$  VI
- Auch die **Streuung** wird auf analoge Weise wie vorher geschätzt  $\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n R_i^2$  **p = Anzahl Koeffizienten  $\beta$**

**Beispiel Sprengungen:** Erschütterung = Distanz und Ladung:  $\log(ersch_i) = \beta_0 + \beta_1 \log(dist)_i + \beta_2 \log(ladung)_i + E_i$

`Spr2.lm2 <- lm(I(Ersch ~ IDist | Ladung, data = SprengS2); summary(Spr2.lm)`

	Min	1Q	Median	3Q	Test auf $H_0: \beta_j = \beta_0$
	-1.00326	-0.26665	0.00522	0.2237	$\beta_0, k \neq j$ beliebige «optimale»
					Werte annehmen können
<b>Coefficients:</b>	[Schätzer $\hat{\beta}_j$ ] [Std Fehler $\hat{\sigma} \hat{\beta}_j$ ] [t value] [P-Wert]				
(Intercept)	$\hat{\beta}_0 = 6.82$	$se(\hat{\beta}_0) = 0.5132$	12.707	<2e-16	***
1Dist	$\hat{\beta}_1 = -1.51$	$se(\hat{\beta}_1) = 0.1111$	-13.592	<2e-16	***
1Ladung	$\hat{\beta}_2 = 0.80$	$se(\hat{\beta}_2) = 0.3042$	2.658	0.0109	*

Schätzung der Standardabweichung der Fehler =  $\sigma \downarrow$  1 ' ' 1  
 Residual standard error: 0.3521 on 45 degrees of freedom  
 Multiple R-squared: 0.8046 [R<sup>2</sup> adjusted R-squared: 0.7962]  
 F-statistic: 92.79 on 2 and 45 DF, p-value: < 2.2e-16  
 F-Test auf «Alle  $\beta_j$  ausser  $\beta_0$  sind gleich 0.»

**Multiple Regression ist nicht Summe einfachen Regressionen!**

- **Nur eine multiple Regression zeigt den wahren Einfluss der einzelnen erklärenden Variablen.**

Vertrauensintervall für  $\beta_j$  - `confint(Spr.lm2, level = 0.95)` #wie bei einfacher Regression

- Verteilung eines einzelnen Koeffizienten  $\beta_j$ :  $\hat{\beta}_j \sim \mathcal{N}(\beta_j, \sigma^2 V_j)$ , wobei Grösse nur von erklärenden Variablen  $V_j$  abhängt.
- Also lautet das  $(1 - \alpha) \cdot 100$  %-Vertrauensintervall:  $\hat{\beta}_j \pm c_w \cdot se(\hat{\beta}_j)$  mit
  - o  $se(\hat{\beta}_j) = \hat{\sigma} \sqrt{V_j}$  o  $c_w$  dem  $(1 - \alpha/2)$ -Quantil der t-Verteilung mit  $(n - p)$  Freiheitsgraden.

**Bestimmtheitsmass – «Multiple R-Squared»**

- misst Anteil durch die Regressionsfunktion erklärten Streuung an der Streuung der Y-Werte.
- Ist wie in einfachen Regression definiert,  $R^2$  liegen zwischen 0 und 1; je grösser, desto besser.
- Es entspricht dem Quadrat der Korrelation zwischen den beobachteten  $y_i$  und den angepassten Werten  $\hat{y}_i$  d.h. es misst den **linearen Zusammenhang**. Die nach **Kleinsten Quadraten** geschätzten Koeffizienten **minimieren** nicht nur die **Quadratsumme der Residuen**, sondern **maximieren** auch die **Korrelation** zwischen den **angepassten Werten** und den **Beobachtungen der Zielgrösse**. Der **maximale Wert** ist die **multiple Korrelation**.

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{SS_{Fit}}{SS_Y}$$

Vielfalt der Modellierungsmöglichkeiten

- Die erklärenden Variablen müssen von **keinem bestimmten Datentyp sein, nicht einer bestimmten Verteilung** folgen (sind keine Zufallsvariablen) und es gibt **keine Voraussetzung** über ihre **Abhängigkeit** untereinander.

**Polynomiale Regression**

- Die  $x$ -Variablen können im Modell in irgendeiner Weise aus ursprünglichen erklärenden Variablen abgeleitet werden.
- Z.B. quadratische Regression:  $Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + E_i$  mit  $x^{(1)} = x$  und  $x^{(2)} = x^2$  ergibt sich
- $Y_i = \beta_0 + \beta_1 x_i^{(1)} + \beta_2 x_i^{(2)} + E_i$  ein multiples lineares Regressionsmodell.
- In gleicher Weise können auch höhere Potenzen eingeführt werden  $\rightarrow$  polynomiale Regression
- Polynomiale Regression ist ein multiples lineares Regressionsmodell. In der polynomialen Regression wird eine **nicht-lineare Funktion** in  $x$  angepasst, die jedoch **linear** in den **unbekannten Parametern  $\beta_0, \beta_1, \beta_2$**  ist!

- **Modell:**  $\log(P_i) = \beta_0 + \beta_1 \cdot \log(C_i) + \beta_2 \cdot \log(C_{2i}) + E_i \Leftrightarrow P_i = e^{\beta_0} \cdot C_i^{(\beta_1 + \beta_2 \cdot \log(C_{2i}))} \cdot e^{E_i}$   $E_i \sim \mathcal{N}(0, \sigma^2)$  unabhängig
- `C$ICarat <- log(C$Carat); C$ICarat2 <- log(C$Carat)^2; C.lm <- lm(I(Price ~ ICarat + ICarat2, data = C))`
- **Linien in Plot:** `x <- seq(min(C$ICarat), max(C$ICarat), length = 50);`  $\rightarrow$  Bei qu. Regr. min-max-seq von  $x$  term
- `predict(C.lm, newdata = data.frame(ICarat = x, ICarat2 = x^2), interval = "c")` #nicht von  $x^2$  bilden und dann erst  $\wedge 2!$

**Erklärende Variable ist binär** `Spreng2$Stelle <- as.factor(Spreng2$Stelle)`

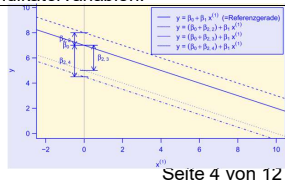
- Einfache lineare Regression mit einer binärer erklärenden Variablen  $x$ :  $Y_i = \beta_0 + \beta_1 x_i + E_i$ ,  $i = 1, 2, \dots, n$  seien die 15 ersten Elemente in  $x$  gleich 0 und die restlichen gleich 1, dann  $Y_i = \beta_0 + \beta_1 \cdot 0 + E_i$ ,  $i = 1, 2, \dots, 15$  und  $Y_i = \beta_0 + \beta_1 \cdot 1 + E_i$ ,  $i = 16, 17, \dots, n$ . Das Regressionsmodell beschreibt zwei unabhängige Stichproben.
- Der Koeffizient  $\beta_1$  beschreibt eine zusätzliche Effekt-Verschiebung in der zweiten Stichprobe.

**Beispiel Sprengungen:** Kann sein, dass örtlichen Gegebenheiten der Messstellen Einfluss auf Erschütterung haben?

- Die Variable, die die Stellen bezeichnet, nennt man **Faktorvariable**. Um sie in ein Regressionsmodell einzubeziehen, führt man für jede mögliche Stelle eine «**Indikatorvariable**» (Dummy Variable) ein:  $x_i^{(j)} = \begin{cases} 1 & \text{falls } i\text{-te Beobachtung aus der } j\text{-ten Stelle (Gruppe)} \\ 0 & \text{sonst} \end{cases}$   $\rightarrow$  Faktorvariable führt zu einem Block von Indikatorvariablen.

**Skizze zur Illustration einer möglichen Parametrisierung von parallelen Geraden**

- Betrachten Fall mit Distanz und Stelle. Modell:  $Y_i = \beta_0 + \beta_1 x_i^{(1)} + \beta_2 x_i^{(2)} + E_i$ , wobei  $Y = \log(\text{Erschütterung})$  und  $x^{(1)} = \log(\text{Distanz})$  und  $x^{(2)}$  die Messstelle bezeichnet
- Sei Stelle 1 die **Referenzstelle** mit  $Y_i = \beta_0 + \beta_1 x_i^{(1)} + E_i$ . Stellen 2, 3 und 4 sind geschifte Geraden, die mit Indikatorvariablen  $St2, St3, St4$  beschrieben werden.
- Der Einfluss der Stelle 1 ist im Achsenabschnitt  $\beta_0$  enthalten.





Das Regressionsmodell mit parallelen Geraden wird geschrieben:

Y\_i = beta\_0 + beta\_1 x\_i^{(1)} + beta\_{2,2} St2\_i + beta\_{2,3} St3\_i + beta\_{2,4} St4\_i + E\_i

beta\_{2,2}, beta\_{2,3}, beta\_{2,4} sind die Verschiebungen der Achsenabschnitten.

Spr.lm <- lm(I(Ersch ~ IDist + ILadung + Stelle, data=Spreg2))

Sind zwei Geraden gleich?

Unterscheiden sich zwei Geraden im Achsenabschnitt/Steigung oder in beidem? Formulieren als Modell:

Y\_i = alpha + beta x\_i + Delta alpha g\_i + Delta beta x\_i g\_i + E\_i

g\_i = 0 falls Gruppe A und g\_i = 1 falls Gruppe B

Gleichung fallweise aufgeschrieben: g\_i = 0: Y\_i = alpha + beta x\_i + E\_i

und g\_i = 1: Y\_i = (alpha + Delta alpha) + (beta + Delta beta) x\_i + E\_i -> Der Term <math>x\_i g\_i</math> wird auch als Wechselwirkung bezeichnet.

Die beiden Geraden stimmen in Steigung überein, wenn Delta beta = 0 oder gesamthaft überein, wenn zugleich Delta alpha = Delta beta = 0

Wechselwirkung kann in R mit lm(Y ~ x1 + x2 + x1 : x2, data = dat) oder lm(Y ~ x1 | x2, data = dat) erfasst werden.

Vergleich von Regressionsmodellen mit F-Test

Messstellen Einfluss auf Erschütterung? Kein Einfluss -> alle Koeffizienten Indikatorvariablen null: Delta\_1 = Delta\_2 = Delta\_3 = Delta\_4 = 0

H\_0: beta\_{j1} = 0, beta\_{j2} = 0, ..., beta\_{jq} = 0, ..., q, ist ungleich 0. Teststatistik F := (SS\_E^\* - SS\_E) / (q \* SS\_E / (n - p))

SS\_E^\* := sum\_{i=1}^n R\_i^2, r\_i aus <math>\llcorner</math>kleinen</math> Modell, (p - q) Koeffizienten | SS\_E := sum\_{i=1}^n R\_i^2, r\_i aus <math>\llcorner</math>grossen</math> Modell, p Koeffizienten.

Verteilung von F unter der Nullhypothese F ~ F\_{q, n-p}, F-Verteilung mit q und n - p Freiheitsgraden.

F-Test mit anova(fit1, fit2) (Modell mit zusätzlicher erkl. Variable vs. Modell ohne diese) durchgeführt. <math>\llcorner</math>ANOVA</math> = Analysis of Variance, Varianzvergleich. Alternative: drop1(fit, test = "F") -> zeilenweiser Test: H\_0 <math>\llcorner</math>beta\_j ist gleich 0</math>

Globaler F-Test - Summary-Output

F-Test im Summary-Output dreht sich um Frage: Beeinflusst die Gesamtheit der erklärenden Variablen die Zielgrösse?

Nullhypothese: H\_0: Alle beta\_j ausser beta\_0 sind gleich 0. Alternative: Mindestens eines dieser beta\_j ist ungleich 0.

Testgr.: F := (SSy - SS\_E) / m, SSy := sum\_{i=1}^n (Y\_i - Y\_bar)^2, SS\_E := sum\_{i=1}^n R\_i^2, Unter H\_0 Testgrösse F -verteilt, m = p - 1 und n - p FG

Theoretische Verteilung der Residuen

Residuen Kleinste-Quadrate-Schätzung sind normalverteilt, haben Erwartungswert 0, aber Varianz ist nicht konstant.

Damit Residuen gleiche Verteilung haben wie Fehler -> skalieren: R\_i := R\_i / sqrt(1 - H\_ii) ~ N(0, sigma^2), i = 1, ..., n

Stand. Residuen: R\_i / sqrt(1 - H\_ii), i = 1, ..., n (sind weder standardnormalverteilt noch t-verteilt). Residuen sind korreliert.

In der Residuenanalyse verwenden wir für Überprüfung der Verteilung die standardisierten Residuen, ausser für das T-A-Diagramm, damit man dort die Abweichungen von einem konstanten Erwartungswert besser erkennen kann.

PRESS-Residuen / Studentisierte Residuen

Residuen, die Vorhersagefehler schätzen (prediction errors). Mit PRESS-Residuen, e\_hat\_i, wird Differenz bezeichnet zwischen Y\_i und angepassten Wert. Ergibt sich, wenn die i-te Beobachtung zum Anpassen des Modells nicht verwendet.

Residuen muss nicht für jede Beobachtung neu angepasst werden: e\_hat\_i := R\_i / (1 - H\_ii), i = 1, 2, ..., n, Varianz: var(e\_hat\_i) = sigma^2 / (1 - H\_ii)

Standardisierung PRESS-Residuum, Ersetzung sigma^2 durch sigma\_hat^2 := e\_hat\_i^2 / (1 - H\_ii) = R\_i^2 / (1 - H\_ii)^2, sigma\_hat^2 := 1 / ((n - p - 1) \* sigma\_hat^2)

Oder, standardisierte/studentisiertes Residuum: R\_i^\* := R\_i / (sigma\_hat \* sqrt(1 - H\_ii)), ist t-verteilt mit (n - p - 1) Freiheitsgraden

Residuen-Analyse bei der multiplen Regression - fit <- lm(I(Ersch ~ IDist + ILadung, data = Spr)

Gleiche Analyse mit bekannten drei Diagrammen. Gibt zusätzliche Diagramme.

Fazit nach Analyse: Welches entscheidendste Unstimmigkeit und wie bereinigen?

Falls sich im Tukey-Anscombe-Diagramm Abweichungen von der angenommenen Form der Regressionsfunktion zeigen, welche erklärenden Variablen sind zu transformieren? -> darum Residuen gegen alle erklärenden Variablen auftragen

Plot erstellen: par(mfrow = c(1, 2)); scatter.smooth(dat\$IDist, resid(fit), lpar = list(col = 2)); abline(h = 0)

Plot pro erkl. Variable erstellen scatter.smooth(dat\$ILadung, resid(fit), lpar = list(col = 2)); abline(h = 0)

Partial Residual Plots (PRP)

Andere erkl. Var. haben Einfluss und können Bild (Grafik violett) verfälschen. Abhilfe PRP. Effekt der anderen erkl. Var. wird herausgerechnet: r\_i^{(-k)} = r\_i + beta\_k x\_i^{(k)}

Falls lin. Regressionsmodell geeignet, streuen Punkte im Diagramm r\_i^{(-k)} gegen

x^{(k)} um Gerade mit Steigung beta\_k. Beispiel für I(Ersch ~ ILadung + IDist ->

Falls notwendig, die erklärende Variable x^{(k)} zu transformieren, so ist die Transformation in dieser Darstellung sichtbar.

par(...); terplot(fit, partial.resid = T, smooth = panel.smooth, ylim = "free", col.res = 1, col.term = 4, col.smth = 2)

Additivität der erklärenden Variablen

Die Effekte von zwei erkl. Variablen addieren sich. Streudiagramm:

(x^{(1)}, x^{(2)}) mit Residuum als strichförmiges Symbol. Länge Strich proportional zum Absolutbetrag, Steigung Strich (+/-1) fürs Vorzeichen

Interaktion mit zusätzlicher Variable modellieren: x\_i^{(3)} := x\_i^{(1)} \* x\_i^{(2)}

library("sfsmisc"); bsp.lm <- lm(y ~ x1 + x2, dat = bsp); p.res.2x(x ~ x1 + x2, data = bsp.lm, scol = 2:1)

Alternative: p.res.2x(x = bsp\$x1, y = bsp\$y2, z = resid(bsp.lm), size = 1.5, slwd = 2) #size/slwd -> klarere Striche

Table with columns: Coefficients, Estimate, Std. Error, t value, Pr(> |t|), P-Vert. Rows include (Intercept), IDist, ILadung, StelleSt2, StelleSt3, StelleSt4.

Schätzung der Standardabweichung der Fehler = sigma -> Residual standard error: 0.3379 on 42 degrees of freedom Multiple R-squared: 0.83 = R^2 Adjusted R-squared: 0.8122 F-statistic: 41.66 on 5 and 42 DF, p-value: 3.194e-15 Globaler F-Test auf H\_0 <math>\llcorner</math>Alle beta\_j ausser beta\_0 sind gleich 0</math> ->

Gewichtete lineare Regression

Was soll man tun, wenn die Residuenanalyse die Notwendigkeit zeigen oder die Art der Datenerhebung impliziert, dass die Varianzen der einzelnen Zufallsfehler als nicht konstant anzunehmen sind? sigma\_i^2 = var(E\_i), i = 1, 2, ..., n

Ansatz: sigma\_i^2 = sigma^2 \* 1/w\_i, wobei w\_i bekannt und sigma^2 aus Daten geschätzt werden. Geht auch als Matrix:

W := matrix(1/w\_i, nrow=n, ncol=1)

Beispiel Fettgehalt in Fischen

Untersuchen, ob sich mittlere Fettgehalt zwischen vier Fischarten unterscheidet. Hierfür werden 3 Fische pro Art zufällig ausgewählt. Von diesen sollen je vier Fleischproben entnommen und mittlere Fettgehalt bestimmt.

Problem: Wegen eines Missgeschicks konnten nicht immer jeweils alle 4 Proben ausgewertet werden.

Offensichtlich: die Proben von einem Fisch sind ähnlicher als die Proben zwischen zwei Fischen. D.h. die Messfehler sind nicht identisch verteilt. Ausweg: Analysiere nur den mittleren Fettgehalt pro Fisch. Fat.mean := 1/No\_i \* sum\_{j=1}^{NoSi} y\_{ij}, wobei y\_{ij} der Fettgehalt des j-ten Fisches aus der Art i ist und NoSi die Anzahl noch Fische der Art i.

Konsequenz: Da der beobachtete mittlere Fettgehalt (Fat.mean) einer Fischart i (Species) auf einer unterschiedlichen Anzahl Proben NoSi (NoS im Datensatz) basiert, haben die Beobachtungen unterschiedliche Varianzen.

Die Varianz des Mittelwerts ist: var(y\_bar\_i) = 1/NoSi \* var(E\_i) = 1/NoSi \* sigma^2, weil die einzelnen Beobachtungen unabhängig sind.

Gewichte festlegen: Mittelwert streut weniger, umso grösser Anzahl Proben NoSi ist. Deshalb Gewicht w\_i = NoSi

Daraus folgt: 2 \* Q(beta) = sum w\_i R\_i^2 = R^T W E bezüglich beta zu minimieren. Schätzung: beta\_hat = (X^T W X)^-1 X^T W Y

entspricht erwartungstreue Regression von sqrt(w\_i) y\_i auf sqrt(w\_i) x\_i^{(1)}, ..., sqrt(w\_i) x\_i^{(m)}, für Varianz gilt var(beta\_hat) = sigma^2 (X^T W X)^-1

Residuen sind: R = (I - H\_W) Y, mit H\_W = X (X^T W X)^-1 X^T

Kovarianzmatrix wird zu var(R) = sigma^2 (W^-1 - H\_W), auf Diagonalen der Matrix stehen die Varianzen der Residuen.

FFat2.lm <- lm(Fat.mean ~ Species, data = FFat2.Agg, weights = NoS); summary(FFat2.lm) #NoS = Var. in df

Welche Residuen soll man in grafischen Darstellungen verwenden? Generell gilt:

für Beurteilung der Verteilung und Streuung der Fehler verwendet man skalierte Residuen: R\_i := R\_i / sqrt(1 - (H\_W)\_ii)

Z.B. im Normalplot, scale-location plot oder im Diagramm Residuen gegen Hebelarme.

Geht es um Eignung der Regressionsfunktion, werden Darstellungen sqrt(w\_i) r\_i gegen sqrt(w\_i) y\_i oder gegen sqrt(w\_i) x\_i^{(k)} betrachtet. In diesen Skalen erfolgt eigentliche Anpassung. Z.B. im T-A-Plot oder im Diagramm Residuen gegen erkl. Variablen.

Weitere Motivation für Gewichte: Residuen-Analyse zeigt Varianzen der Zufallsfehler sigma^2 = var(E\_i) nicht konstant ist.

Dann ist allerdings die Herausforderung, die relativen Genauigkeiten v\_i für jede Beobachtung i zu finden, sodass var(E\_i) = sigma^2 v\_i bestimmt werden kann.

Falls sich in einem Streudiagramm der Residuen gegen ein X^{(j)} zeigt, dass die Streuung von X^{(j)} abhängt, dann kann man versuchen, eine Funktion v(X^{(j)}) anzugeben, die diese Abhängigkeit beschreibt.

dann Anwendung der gewichteten Regression mit Gewichten: w\_i = 1/v(x\_i^{(j)})

o sigma^2 ist dann die Varianz jener Beobachtung, bei der v\_i = 1 ist.

Beispiel: Hängt die Varianz von einer erklärenden Variablen ab? scatter.smooth(SyD2\$x1, sqrt(abs(rstandard(SyD2.lm))), lpar = list(col = "red"))

scatter.smooth(SyD2\$x2, sqrt(abs(rstandard(SyD2.lm))), lpar = list(col = "red"))

Ja, die Streuung hängt von der erkl. Var. x1 ab.

d.h. mit steigenden x1-Werten nimmt die Streuung proportional zu

Problemen deshalb Ansatz: SyD2.lm2 <- lm(y ~ x1 + x2, weight = 1/SyD2\$x1^2, data = SyD2)

Gewichtete Regression führt zu besserem Modell als ohne Gewichte. (Kontrolle mit den drei bekannten Plots)

Beispiel: Produktion mit drei Geschwindigkeiten: 100, 150, 200. Wie viele Defekte gibt es? #Var nicht konstant in Scale-P Modell: g.lm <- lm(DEFECTS ~ SPEED, data = G); (G.var <- aggregate(resid(g.lm) ~ G\$SPEED, FUN = var)); #Ausproben was hilft (Werte ca. gleich) -> G.var[,2]/c(100,150,200)^2; G.var[,2]/(c(100,150,200)/100)^2 #passt am besten (G\$w <- 1/((G\$SPEED/100)^2)); g.lmw <- lm(DEFECTS ~ SPEED, data = G, weights = w) #Modell passt nun

Einflussreiche Beobachtungen

Thema Ausreisser: Die Antwort auf die Frage, ob eine Beobachtung ein Ausreisser sei, hängt vom Modell ab!

Wie stark beeinflusst eine Beobachtung (Ausreisser) Analyse? -> Analyse ohne die fragliche Beobachtung wiederholen und Effekte auf Schätzung, Tests, Vertrauensintervalle, ..., untersuchen

Effekt messen mit Cook's Distanz: misst Veränderung aller angepassten Werte y\_hat\_i beim Weglassen i-ten Beobachtung. d\_i = 1/p \* sigma\_hat^2 \* sum\_{j=1}^p (y\_hat\_j^{(-i)} - y\_hat\_j)^2, i = 1, ..., n, wobei y\_hat\_j^{(-i)} angepasster Wert bei Weglassen i-ten Beobachtung ist

Faustregel: Zu einflussreiche Beobachtung = Cook's Distanz grösser als 1 d\_i = 1/p \* R\_ii^2 / (1 - H\_ii), mit R\_ii = R\_i / sqrt(1 - H\_ii)

Hebelarm H\_ii: Diese Grösse misst, wie <math>\llcorner</math>untypisch</math> die Beobachtung in Bezug auf die erklärenden Variablen ist.

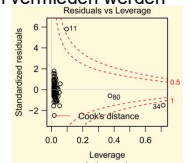
Wertebereich von H\_ii: 0 <math>\llcorner</math>H\_ii <math>\llcorner</math>1, Varianz der Residuen: var(R\_i) = (1 - H\_ii) sigma^2 -> zu einflussreiche Beobachtung?

Wenn ein Wert Y\_i um Delta y\_i verändert, dann misst H\_ii Delta y\_i, die Veränderung des zugehörigen angepassten Wertes y\_hat\_i.

Faustregel I: Beobachtungen mit einem Hebelarm H\_ii > 2^p/2^n haben zu grosse Hebelwirkung plot(fit, which = 4)

Faustregel von Huber: Beob. mit Hebelarm unter 0.2 sind unbedenklich; jene mit über 0.5 sollten vermieden werden Achtung: Hebelpunkte H\_ii und Cook's Distanz untersuchen Effekt nur einer Beobachtung!

Beispiel: Es sind zwei zu einflussreiche (Cook's Distanz > 1) Beobachtungen sichtbar: Beobachtung 11 ist wegen grobem Ausreisser (in y-Richtung) zu einflussreich, Beobachtung 34 ist wegen zu grossen Hebelarm (h\_i > 0.2) zu einflussreich, Beobachtung 80, ist ein weiterer, jedoch ungefährlicher Hebelpunkt.



Robuste Anpassungsmethoden

- In **normaler** Residuen Analyse werden **eventuelle Ausreiser** und **grosse Hebelpunkte gefunden und entfernt**.
- **Robuste Methode:** findet **Ausreiser leichter** und können diese aufgrund der Robustheit der Schätzer **stehen lassen**.
- **Ausreisser** stellen nur eine «**kleine**» **Abweichung** von den Annahmen dar, aber mit **verheerender Wirkung**.
- Oft existiert keine vernünftige Alternative zur Normalverteilung oder zur Linearität. Deshalb geht man davon aus, dass
  - o das **Regressionsmodell** für die **Mehrheit der Beobachtungen stimmt**,
  - o es aber **möglich** ist, dass in **wenigen Fällen Abweichungen** davon **auftreten können**.
- Gesucht ist Schätzverfahren, die Parameter so schätzt, wie wenn Abweichungen nicht vorhanden sind -> robust
- Erster informeller Ansatz zur Robustheit ist, **erstens die Daten auf offensichtliche Ausreisser** zu prüfen, **zweitens** diese zu **löschen** und **drittens** das zum **Modell optimale Schätzverfahren** auf die bereinigten Daten anzuwenden.
- Dieses Verfahren ist allemal besser, als die Ausreisser zu ignorieren. Die Resultate sind jedoch immer zu optimistisch.
- Ob Schätzer **robust** ist, kann mit **zwei Massen** untersuchen. Beide Masse beruhen auf Idee, **Verhalten einer Schätzfunktion** unter dem **Einfluss von groben Fehlern** (gross errors), d.h. willkürlichen hinzugefügten Daten, zu studieren:
  - o **Einflussfunktion** (influence function): Die Sensitivität (gross error sensitivity) basiert auf der Einflussfunktion und misst den maximalen Effekt einer einzelnen Beobachtung auf den Schätzwert.
  - o **Bruchpunkt** (breakdown point): Der Bruchpunkt gibt den **minimalen Bruchteil von Beobachtungen** an, der genügt, um **unglaubliche Schätzwerte** zu erhalten. Bruchpunkt misst **Grösse der Störung**, die zur **Katastrophe** führt.
- **Schätzer guten Robustheitseigenschaften:** beschränkte Sensitivität und Bruchpunkt möglichst nahe an max. Wert  $\frac{1}{2}$ .

Konstruktion eines robusten Regressionschätzers

- **Idee:** Gewicht der Ausreiser mittel gewichteter linearer Regression herunterzugewogen respektive beschränken.
- Gewichtete Normalgleichung  $X^T W R = 0$  respektive  $\sum_{i=1}^n w_i \cdot R_i \cdot x_i$ , mit  $R_i = Y_i - x_i^T \underline{\beta}$  zeigt, dass Beobachtungen  $i$  mit grossen Residuen  $R_i$  mit entsprechenden Gewichten  $w_i$  heruntergewichtet werden.
- Geeignete Gewichte mittels Einflussfunktion  $IF$  bestimmen  $IF(y) = w_i \cdot r_i \cdot M \cdot x_i$ , wobei  $M$  noch zu bestimmende Matrix.
- Wollen Einfluss grosser Residuen beschränken, muss  $\psi(r_i) := w_i \cdot r_i$  beschränkt werden. Gewicht mit  $w_i = \frac{\psi(r_i)}{r_i}$  bestimmen. Da Gewicht von Residuen abhängt und damit vom wahren Parameter  $\underline{\beta}$  mit  $r_i = y_i - x_i \cdot \underline{\beta}$ , können wir Gewicht nicht mehr explizit berechnen. Das führt dazu, dass Lösung nur mit iterativen Algorithmus gefunden kann und Standard-Fehler der Schätzungen anders berechnet werden müssen als in der gewichteten linearen Regression.

Robuster Regressions-M-Schätzer (ist Schätzer mit oben beschriebenen impliziten Gewichten)

- In Funktion  $\psi$  gibt es einen Knick, welcher festlegt, von welcher Stelle extreme Beobachtungen an Einfluss verlieren.
- Sinnvoll wenn Fehler in ihrer Einheit  $\frac{r_i}{\sigma}$  gemessen werden. Knickstelle  $c$  gegeben durch theoret. Effizienz  $c = 1.345$
- M-Schätzgleichung:  $\sum_{i=1}^n \psi(\frac{\hat{\beta}_M}{\sigma}) \cdot x_i = 0$ , mit  $\hat{r}_i = \frac{Y_i - x_i^T \hat{\beta}_M}{\sigma}$  oder mit Gewichten:  $\sum_{i=1}^n w_i \cdot r_i \cdot x_i = 0$ , mit  $w_i = \frac{\psi(\frac{\hat{\beta}_M}{\sigma})}{\hat{r}_i}$
- Die Standardabweichung als Schätzung des Skalenparameters (Streuung)  $\sigma$  ist jedoch extrem unrobust. Eine bessere robustere Schätzung für den Skalenparameter ist zum Beispiel der **median of absolute values**:  $\hat{\sigma}_{MAV} := \frac{med(|r_i|)}{0.6745}$
- Wenn sich ein Ausreisser auf Hebelpunkt befindet, erweist sich die Regressions-M-Schätzung als **ungenügend robust**.
- Dieser Schätzer begrenzt nur Einfluss Residuen, nicht Einflussfunktion als Ganzes, weshalb sie Bruchpunkt von 0 hat.

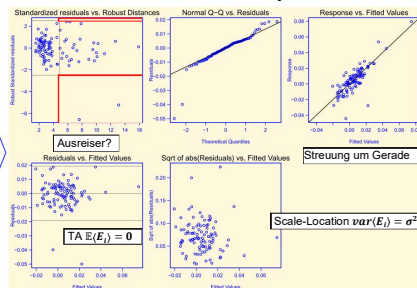
library("MASS"); f.rlm <- rlm(Y ~ X1 + X2, method = "M", data = dat); summary(f.rlm); par(mfrow = c(2, 2)); plot(f.rlm)

Regressions-MM-Schätzung (Regressions-Modified-M-Estimator)

- Um Schätzer zu konstruieren, der robust mit hohen Bruchpunkt ist, müssen Beobachtungen mit grossen Residuen nicht nur heruntergewichtet, sondern allmählich (in Abhängigkeit des Betrages der Residuen) ignoriert werden.
- Da  $\psi$ -Funktion wieder auf die Nulllinie absinkt, werden solche Regressions-M-Schätzern als **re-descending** bezeichnet.
- Da Schätzfunktion (viele) lokale Minima besitzt, ist Optimierung schwierig. Lösung Problem: mit etwas anderen Optimierungsaufgabe mit stochastischem Suchalgorithmus werden gute Startwerte für Koeffizienten  $\underline{\beta}$  gefunden.
- Die Lösung ist als S-Schätzer (ist Startwert) bekannt, hat jedoch schlechte theoretische Eigenschaften.
- Die zum Startwert nächst gelegene lokale Lösung des eigentlichen Problems führt dann zu einem besseren Schätzer, wobei als Schätzwert für  $\sigma$  jene aus den Startwerten, d.h.  $\hat{\sigma}_0$  verwendet wird.
- Lösung hat hohen Bruchpunkt von  $\epsilon^* = \frac{1}{2}$  und die Inferenz des Schätzers kann analog zu M-Schätzern gemacht werden.
- Chance, aus allfälligen Abweichungen besseres Modell entwickeln. Da robuste Anpassungsmethoden durch «Verunreinigungen» der Daten kaum irritiert werden, sind Modellabweichungen viel leichter in der Residuen-Analyse zu erkennen.
- Schätzer approximativ normalverteilt:  $(\hat{\beta}_M)_j \sim \mathcal{N}(\beta_j, \sigma^2 \tau_j)$ , mit  $\tau > 1$  Korrekturfaktor. Approximation besser, je grösser  $n$
- Vertrauensintervall  $(\hat{\beta}_M)_j \pm c_w \cdot se((\hat{\beta}_M)_j)$  mit  $se((\hat{\beta}_M)_j) = \hat{\sigma}_0 \sqrt{\tau_j}$ ,  $c_w$  Quantil t-Verteilung mit  $n - p$  Freiheitsgraden und  $\hat{\sigma}_0$  der Residuenstandardfehler vom initialen S-Schätzer ist.
- Robuste Version des F-Tests ist eine  $D(x_j, \hat{\beta}_{MM}) := 2 \cdot \hat{\sigma}_0^2 \cdot \sum_{i=1}^n \rho(\frac{y_i - x_i^T \hat{\beta}_{MM}}{\hat{\sigma}_0})$   $\Delta^2 \sim \chi^2_p$  Verteilung unter der  $H_0$

library("robustbase"); #KS2014 für effizientes Vertrauensintervall f.lmrob <- lmrob(Y ~ X1 + X2, setting = "KS2014", data = dat) par(mfrow = c(2,3)); plot(f.lmrob, las = 1)

- > Zweck Robuste Plots: Identifizierung der Ausreiser
- > sehen hier im Plot überall zwei Ausreiser



Variablenselektion und Modellbildung

- Welche erklärenden Variablen sollen wie ins Regressionsmodell eingehen? Warum kein Modell mit allen erklä. Var.?
- o **Einfachheit** – gibt es mehrere Erklärungen, so sollte man die Einfachste wählen.
- o **Reduktion Schätzvariabilität** – unnötige erklärende Variablen verschlechtern Genauigkeit der Modellschätzung.
- o **Bessere Interpretierbarkeit** – Multikollinearität führt zu nicht eindeutigen Lösungen.
- o **Vorhersage** – weniger erklärende Variablen, weniger Aufwand beim Vor- und Aufbereiten der Daten.
- Ist ein bestimmter Term  $\beta_j x^{(j)}$  im Modell nötig? Nützlich? Überflüssig? -> Grundbaustein für die Variablenselektion
- Als Hypothesen-Prüfung diese Frage schon gelöst: Nullhypothese  $\beta_j = 0$  prüft (t-Test) -> multiples Testproblem
- Weiteres Problem: müsste Voraussetzungen über Fehler prüfen, wenn P-Werte der Tests zum Nennwert nehmen wollte.
- **Lösung:** Statt Tests für strikte statistische Schlüsse zu verwenden, nehmen P-Werte der t-Tests.
- Bei Faktorvariablen Ersetzung t-Test durch F-Test: **add1(..., scope=..., test = "F")** oder **drop1(..., test = "F")**

Manuelle Vorwärts-Selektion (Forward Selection)

- Mit diesem auf P-Werten gestützten Kriterium können wir eines der Variablenselektionsverfahren wie folgt formulieren:
  - Wir starten mit dem einfachsten kleinen Modell:  $Y_i = \beta_0 + E_i$  **AKW.lm0 <- lm(Y ~ 1, ...)**
  - Ersten Schritt wird die Variable in das Modell aufgenommen, die in einfacher Regression den **kleinsten P-Wert** besitzt.
  - Dann wird jeweils die nächste noch nicht berücksichtigte Variable mit dem kleinsten P-Wert aufgenommen.
  - **Stoppregel:** Grundsätzlich nicht festlegbar. Schranke 0.05 pro Zeile ist üblich, problematisch da multiples Testproblem
- AKW.lm0 <- lm(lgK ~ 1, data = AKW)** #Beginn mit Zielvariable ~ 1, einfaches Modell erstellen
- add1(AKW.lm0, scope = ~ lgG + D + WZ + BZ + Z + NE + KT + BW + sqrtN + KG, test = "F")** #add1 Fun wiederholen
- AKW.lm0 <- update(AKW.lm0, ~. + KG); drop1(AKW.lm0, test = "F")** #Hinzufügen Variable + Überprüfung p-Wert
- add1(... genau wie oben); AKW.lm0 <- update(AKW.lm0, ~. + lgG); drop1(AKW.lm0, test = "F")** # p-Werte i.O.?
- ...add1(... genau wie oben); AKW.lm0 <- update(AKW.lm0, ~. + KT) drop1(AKW.lm0, test = "F")** # p-Werte i.O.?
- AKW.lm0 <- update(AKW.lm0, ~. + KT)** #KT Variable nun nach drop1() nicht mehr signifikant, deshalb wieder entfernen
- summary(AKW.lm0)** -> gibt nun das Modell an, AKW.lm0 ist nun folgendes Modell:  $lgK \sim KG + lgG + D + NE$

Manuelle Rückwärts-Selektion (Backward Selection)

- Zweiter Vorschlag für ein Variablenselektionsverfahren, das auf den P-Werten gestützt ist, lässt sich wie folgt formulieren
  - Wir starten mit dem «vollen» Modell aller erklärenden Variablen  $Y_i = \sum_{k=0}^m \beta_k x_i^{(k)} + E_i$
  - In den weiteren Schritten wird jeweils **eine Variable aus dem Modell genommen**. Es ist jeweils die erklärende Variable, die am Variablenselektion mit P-Werten «unwichtigsten» ist und damit den **grössten P-Wert** hat.
  - **Stoppregel:** Grundsätzlich wieder nicht festlegbar. Oft sieht man die Schranke 0.05 (**P-Wert für jede Zeile**)

- AKW.lm0 <- lm(lgK ~ lgG + D + NE + KG + BW + BZ + Z + sqrtN + KT, data = AKW)** -> zuerst volles Modell erstellen
- drop1(AKW.lm0, test = "F"); AKW.lm1 <- update(AKW.lm0, ~. - BZ)** #Variable mit höchstem P-Wert entfernen!
- ...drop1(AKW.lm0, test = "F"); AKW.lm0 <- update(AKW.lm0, ~. - KT)** #Variable mit höchstem P-Wert entfernen!
- drop1(AKW.lm0, test = "F")** -> nun alle Werte signifikant. AKW.lm0 ist nun folgendes Modell  $lgK \sim lgG + D + NE + KG$
- I.A. **nicht zwingend**, dass beide Selektionsverfahren (Vorwärts- und Rückwärtsselektion) zum **gleichen Modell** führen.
- Achtung, nur eine Variable aufs Mal ausschliessen! Verbleibenden Prädiktoren haben kleinere p-Werte als volles Modell.
- Die Vorhersagekraft konzentriert sich von ursprünglich vielen auf die nun selektierten Variablen.
- Für Vorhersagen ergeben sich oft zu kleine Modelle, die nicht optimale Prognosen geben.
- Entfernte Variablen können kausalen, jedoch kleinen, Einfluss auf die Zielgrösse haben, der im Rauschen untergeht.

Schrittweise Selektion – Kombination aus Vorwärts- und Rückwärtselimination

- Wird in jedem Schritt ausprobiert, ob Weglassen oder Hinzunehmen einer erklä. Var. Modellwahlkriterium verbessert.
- Jene Aktion, die zur grössten Verbesserung im gewählten Modellwahlkriterium führt, wird durchgeführt.
- Verfahren wird beendet, sobald keine Verbesserung mehr möglich ist.

Modellwahlkriterien

- **Variablenselektionsverfahren** beruhen auf **Modellwahlkriterien**, die **Güte der Modellanpassung** durch statistische **Masszahl** beschreibt. Die Masszahl sollte **Modellgenauigkeit** und die **Modellkomplexität** bewerten.
- In der linearen Regression mit dem KQ-Anpassungskriterium kann
  - Modellgenauigkeit mit Summe der **Residuenquadrate**
  - Modellkomplexität mit **Anz. Koeffizienten** gemessen werden
- Ziele Modellkomplexität und Modellgenauigkeit sind konträr. Zusatz. erklä. Var. führen zu höheren Modellgenauigkeit.**

Korrigiertes Bestimmtheitsmass (adjusted R^2<sub>adj</sub>):

- Reguläres  $R^2$  wird immer grösser je mehr Variablen -> ungeeignet da keine Berücksichtigung der Modellkomplexität.
- Deshalb **korrigiertes  $R^2$** :  $n =$  Anzahl Beobachtungen,  $R^2_{adj} = 1 - \frac{SS_E}{(n-1)} = 1 - \frac{n-1}{n-p} (1 - R^2) = 1 - \frac{\hat{\sigma}^2}{\frac{SS_Y}{(n-1)}}$ , wobei  $\frac{SS_E}{(n-p)}$
- $p =$  Anzahl zu schätzenden Koeffizienten  $\beta_k$

Mallow's  $c_p$ -Statistik:

- minimiert den Vorhersagefehler  $C_p := \frac{SS_E}{\hat{\sigma}_p^2} + 2p - n = (n-p) \left( \frac{\hat{\sigma}_p^2}{\hat{\sigma}_p^2} - 1 \right) + p$
- $\hat{\sigma}_p^2 = \frac{SS_E}{(n-p)}$  = Schätzung  $\hat{\sigma}^2$  aus Modell mit  $p$  Parametern,  $\hat{\sigma}_p^2 =$  Schätzung  $\hat{\sigma}^2$  aus grösstem Modell mit  $p^*$  Parametern

Das Informations-Kriterium von Akaike (AIC):  $AIC = -2(\max \text{LogLikelihood}) + 2(\text{Anz. geschätzter Parameter})$

- $AIC = n \cdot \log(\frac{1}{n} SS_E) + 2p^* + \text{Konstante}$ ,  $p^* =$  Anz. gesch. Parameter inkl.  $\hat{\sigma}$ ,  $\diamond =$  Diamant. AIC ist gut verallgemeinerbar

In R können diese drei Variablenselektionsverfahren basierend auf dem AIC mit step(...) durchgeführt werden:

- **Vorwärts:** **M0 <- lm(Y~1, data=DF); step(M0, scope=list(lower=~1, upper=~ x1+x2+x3+x4), direction = "forward")**
- **Rückwärts:** **MF <- lm(Y ~ x1 + x2 + x3 + x4, data = DF); step(MF, direction = "backward")**
- **Schrittweise:** **Mm <- lm(Y ~ x2, data = DF); step(Mm, scope = list(lower = ~ 1, upper = ~ x1 + x2 + x3 + x4))**
- Letztes Modell in Ausgabe ist jeweils das beste Modell.  $\uparrow$  scope = range zwisch. kleinsten/grössten Modell

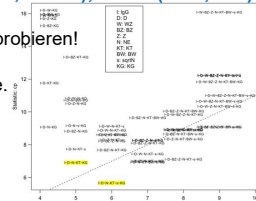


**Eindeutige Lösung? Optimales Modell?**

- Häufig führen verschiedenen Selektionsarten zu unterschiedlichen Ergebnissen.
- Grund: Keines dieser Variablenselektionsverfahren garantiert, dass gemäss Modellwahlkriterium **global optimale Modell** gefunden wird. → Können auf dem Weg in **lokalen Optimum festhängen** bleiben.
- Einzige Möglichkeit global optimiertes Modell finden, ist alle Modellvarianten zu berechnen. → **All-Subset Selection**
- Bei **Vollsuche** ist der Rechenaufwand zu gross → kombinatorische Explosion,  $2^m = \text{Anz. Modelle}$ ,  $m = \text{Anz. erkl. Var.}$

**Vollsuche (All-Subset Selection) mit  $C_p$ -Statistik**

- **Vorteil:** Falls Modell alle notwend. Var. enthält, dann ist erwartete Kriteriumswert  $\mathbb{E}(C_p) = \text{Anz. Parameter } p$ ,  $\mathbb{E}(C_p) = p$
- Fehlen notwendige Variablen, so ist  $\mathbb{E}(C_p) \gg p$  `library(leaps); library(car) #Plot Streudiagramm  $C_p$  gegen  $p$`
- `AKW.Cp <- regsubsets(lgK ~ lgG + D + WZ + BZ + Z + NE + KT+BW + sqrt(N) + KG, nbest=6, nvmax=10, data=AKW)`
- `h <- subsets(AKW.Cp, statistic = "cp", legend = "top", min.size = 4, cex.subsets = 0.7, las = 1); abline(a = 1, b=1)`
- Oder `par(mai = c(1, 1, 1, 3))` und `legend = "interactive"`, falls Legende Kombi überdeckt
- `nbest:` Anzahl Modell mit  $p$  Koeffizienten, die berücksichtigt werden (default = 1) → ausprobieren!
- `nvmax:` Modelle mit maximal  $p = \text{nvmax}$  werden berücksichtigt (default = 8)
- `subsets()` zählt Achsenabschnitt **nicht**, deshalb `abline(a = 1, b = 1)` = Winkelhalbierende.
- `statistic` hat noch Argumente "bic" (Bayes Information Criterion) und "adjr2" ( $R_{adj}^2$ )
- **In Frage** kommenden Modelle müssen um **Winkelhalbierende  $C_p = p$  streuen** und **min.  $C_p$ -Wert** und **min.  $p$  Wert** haben. `#legend = "interactive"` für Plot



**Anmerkung zur Variablenauswahl**

Falls nicht Vollsuche verwendet werden kann:

- Schrittweise Variablenauswahl mit möglichst **grossem Modell** als **Start**, da der Suchprozess etwas umfassender ist.
- Vorwärtsselektion vor allem bei Datensätzen mit **vielen erklärenden Variablen** und **wenig Beobachtungen**.
- Bei Modellen mit **Faktorvariablen**, **Polynom** oder **Wechselwirkungstermen** sind folgende Regeln zu beachten:
- Bei **Faktorvariablen** werden nicht einzelne Faktorstufen **entfernt**, sondern nur die **ganze Variable**.
- **Wechselwirkungsterme** immer mit **zugehörigen Haupteffekte** im Modell **behalten**.
- Polynomen immer alle Terme bis max. Ordnung im Modell **behalten**, d.h. Var.selection bestimmt nur max. Ordnung.
- **Vorsicht** bei **manueller Variablenauswahl**. `step()`, `drop1()`, `add1()` befolgen diese Regeln, **nicht aber regsubsets()**.

**Alternative zu AIC: Bayes Informationskriterium (BIC)**

- $BIC = n \cdot \log(\frac{1}{n} SS_E) + \log(n) \cdot p + \text{Konstante}$ ,  $p = \text{Anz. gesch. Parameter inkl. } \delta, \diamond = \text{Diamant}$ .
- `step(..., k = log(nrow(datasets)))` #Rest ist identisch wie bei AIC, Angabe von  $k$  führt zu BIC

**AIC oder BIC? – Beide führen zu ähnlichen Entscheidungen.**

- BIC, wenn verstehen will, welche Prädiktoren Beitrag leisten und wenn schlankes, gut interpretierbares Modell will.
- AIC, wenn Modell für Vorhersage von zukünftigen Werten will, die verwendeten Prädiktoren aber weniger zentral sind.

**Welches Modell ist das richtige?** Das «beste» Modell ist noch lange nicht das «richtige» oder «wahre» Modell!

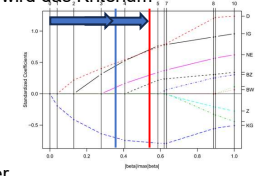
- Deshalb immer mehrere Modelle in Betracht ziehen, die Kriterien «gut» – nicht viel schlechter als «beste» – bewerten.
- Wie viel schlechter? – Leider gibt Theorie darauf keine Antwort. Auch berücksichtigen ist **Kollinearität** unter erkl. Var.

**Alternative Ansätze für Variablenauswahl**

- Heute gibt es sehr viele Daten, auch hochdimensionale Daten und mehr Variablen als Beobachtungen im Datensatz.
- Solche Daten können mit Regressionsmodell analysiert werden, da oft nur wenige erkl. Var. Einfluss auf Zielvar. haben.
- Technik, für solche Situationen, finden sich unter Begriff **«Shrinkage Schätzung»** oder **Regularisierungsmethode**.
- Sie stellt sicher, dass es nicht zu **«Overfitting»** kommt (d.h. das Regressionsmodell wird nicht zu gut an die Daten angepasst, sodass es das Rauschen als systematischen Effekt interpretiert und zu unsicheren Vorhersagen führt)

**LASSO-Methode – «Least Absolute Shrinkage and Selection Operators»**

- ist eine Vertreterin Shrinkage-Schätzer, die zusätzlich **implizit** noch die Eigenschaft hat, **Variablen zu selektionieren**.
- Shrinkage-Schätzer **shrumpft Koeffizienten in Richtung gegen 0**. Wenn ihr **Beitrag zur Erklärung der Zielvariable klein** ist, wird beim LASSO der entsprechende Koeffizient auf 0 gesetzt und hat so **Eigenschaft Variablenauswahl**.
- Kriterium KQ wird durch «Bestrafungsterm» für Grösse der Koeffizienten ergänzt. Somit wird das Kriterium  $Q(\beta; \lambda) := \sum_{i=1}^n \hat{R}_i^2 + \lambda \sum_{k=1}^m |\beta_k|$  mit  $\hat{R}_i := \hat{y}_i - \hat{x}_i^T \beta$  **minimiert**, Zielvariable ist zentriert:



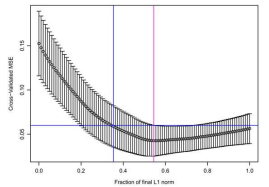
- $\hat{y}_i := y_i - \bar{y}$ , Vektor  $\hat{x}_i$  aus standardisierten erkl. Var. besteht:  $\hat{x}_i^{(k)} := \frac{y_i^{(k)} - \bar{y}^{(k)}}{sd(x_i^{(k)})}$  und Vektor  $\hat{x}_i$  keine 1 als erstes Element hat (Achsenabschnitt) `par(las = 1, cex.axis = 0.5)`
- `library(lars); h.x <- model.matrix(~ -1 + lgG + D + WZ + BZ + Z + NE + KT + BW + sqrt(N) + KG, data = AKW); AKW.lasso <- lars(x = h.x, y = AKW$lgK); plot(AKW.lasso); legend("topleft", legend = AKW.lasso$actions)`
- zwei optimale Parameter  $\tilde{b}$  von Plot auf der nächste Seite (gleiche x-Achsenkala, hier ca. 0.35 und 0.55) hier auf Grafik legen. Dann erkl. Variable bis zum Strich von **links** nach rechts geben das Modell.

**Optimalen Bestimmung des tuning Parameters  $\tilde{b}$**   $\triangleright$  Modell: (lgK-KG + D + IG + NE) und (lgK-KG + D + IG + NE + KT)

- üblicherweise mittels fünf- bis zehnfachen Kreuzvalidierungsverfahren. Das heisst bei fünffachen Kreuzvalidierung, o Fünftel der Daten weggelassen und auf restlichen Daten Lasso-Verfahren für ein Raster von  $\tilde{b}$ -Werten durchführen.
- Anschliessend wird für jede Lasso-Lösung, indiziert über die verschiedenen  $\tilde{b}$ -Werte, Prognosen für die ausgelassenen Beobachtungen gemacht und mit den Zielwerten verglichen.
- Dies wird dann für die anderen 4 Pakete von einem Fünftel der Beobachtungen wiederholt.
- Daraus kann man dann den mittleren Prognosefehler und eine Schätzung dessen Streuung bestimmen.

**Das optimale  $\tilde{b}$  ist dort, wo**

- entweder der **mittlere Prognosefehler minimal** ist (rote Linie in Grafik)
- oder der **mittlere Prognosefehler kleiner** als der **minimale Prognosefehler plus** eine **Standardabweichung** ist. (blaue Linien in Grafik) #Befehle oben gehören auch dazu
- `library(lars); set.seed(4567); AKW.lasso.cv <- cv.lars(x = h.x, y = AKW$lgK, K = 10)`
- `(h.wMin <- which.min(AKW.lasso.cv$cv)); (h.sel <- which(AKW.lasso.cv$cv <= (AKW.lasso.cv$cv[h.wMin] + AKW.lasso.cv$cv.error[h.wMin]))[1])`
- `abline(v = AKW.lasso.cv$index[c(h.wMin, h.sel)], col = c(2, 3))`
- `abline(h = AKW.lasso.cv$cv[h.wMin] + AKW.lasso.cv$cv.error[h.wMin], col = 3)`
- `c(AKW.lasso.cv$index[h.sel], AKW.lasso.cv$cv[h.sel])` #Erster Wert ist das optimale b
- `plot(AKW.lasso, las = 1); abline(v = AKW.lasso.cv$index[c(h.wMin, h.sel)], col = 4, lwd = 5)`
- **Koeffizientenschätzung an Stelle optimalen b:** `coef(AKW.lasso, s=AKW.lasso.cv$index[h.sel], mode="fraction")`



**Kollinearität**

- **Hohe Korrelationen** zwischen erkl. Var. oder allgemeinere Formen davon (Kollinearität) sind von **Theorie zugelassen!**
- Führt aber zu Problemen bei Interpretation und beim statistischen Modellieren. Im Vorwärts- und Rückwärts-Verfahren kann es allenfalls vom Zufall abhängen, welche der beteiligten Variablen als erste weggelassen/aufgenommen wird.
- Bei Kollinearität gilt: Eine erkl. Var. lässt sich annähernd als Linearkombination der anderen darstellen,  $x_i^{(j)} \approx \gamma_0 + \sum_{k \neq j} \gamma_k x_k^{(k)}$ , **gilt die Beziehung exakt**, dann gibt es **keine eindeutige Lösung** bei der KQ-Schätzung.

**Variance inflation factor (VIF)**

- Ist die Beziehung  $x_i^{(j)} \approx \gamma_0 + \sum_{k \neq j} \gamma_k x_k^{(k)}$  nur annähernd erfüllt, können wir das auch als Regressionsproblem ansehen.
- Das Bestimmtheitsmass der Regression von  $x^{(j)}$  auf alle übrigen erkl. Var. ( $R_j^2$  genannt) zeigt,
  - o **wie stark** eine solche **Beziehung** ist und
  - o ist sinnvolles **Mass für Kollinearität**,
  - o das erst noch angibt, **welche Variable** «das **Problem** verursacht».
- Ein daraus abgeleitetes Mass für Kollinearität ist der  $VIF_j = \frac{1}{1 - R_j^2}$

• **Faustregel:** Falls diese Grösse grösser als **5 bis 10** ist, dann sind Probleme mit der Kollinearität vorhanden.

• `library(car); round(vif(AKW.lmV), 2); pairs(AKW); #Streudiagramm oder library(ellipse); plotcorr(cor(AKW))`

**Weitere Masse zur Identifizierung von Kollinearität**

- Die **Konditionszahl (\*)** eignet sich für **Entdeckung von Multikollinearitäten**. Sie ist definiert als  $\kappa = \frac{\lambda_{\max}}{\lambda_{\min}}$ , wobei  $\lambda_{\max}$  der maximale und  $\lambda_{\min}$  der minimale Eigenwert von  $X^T X$  sind. Da  $X^T X$  symmetrisch ist, sind alle Eigenwerte reell.
- Liegt die Konditionszahl **zwischen 100 und 1'000**, so liegt eine **moderate bis starke Multikollinearität** vor.
- Ist sie **grösser** als 1'000, so ist die **Multikollinearität schwerwiegend**.

**Auswirkungen der Kollinearität**

- **Hohe Kollinearität** führt zu **grossen Standardfehlern** bei den **geschätzten Koeffizienten**. Man kann zeigen, dass die Standardfehler für die jeweiligen Koeffizienten  $\beta_j$  mit  $\sqrt{VIF_j}$  aufgeblasen werden. (Optimal wäre ein Faktor von 1.)
- Oft sind deshalb viele/alle Variablen gem. t-Test nicht signifikant. Können jedoch gem. F-Test nicht alle Var. weglassen.
- In gewisse Richtungen sehr **kleine Prognosefehler**, in **andere** sehr **grosse**.
- **Effekte** der einzelnen **Variablen** auf die **Zielvariablen** können **nicht richtig interpretiert** werden.
- **Kollinearität muss bereinigt werden**, damit die **Effekte** und ihre statistische Signifikanz **interpretierbar** werden.

**Was tun gegen Kollinearität?**

- Wenn immer möglich, soll man Beobachtungen so durchführen, dass das Problem vermieden wird. Ansonsten:
  - o Variablen linear transformieren; d.h. z.B. stark korrelierte Variablen **ersetzt** man z.B. durch ihre **Summe** und ihre **Differenz**, falls sinnvoll (u.a. gleiche Einheiten), es können auch andere (lineare) Funktionen verwendet werden:
  - o **Weitere Möglichkeiten:** **Mittelwert** ( $0.5 * (...)$ ) oder **Verhältnis** (Relation) bilden.
  - o `Jet$mean <- 0.5*(Jet$X1 + Jet$X2); Jet$dif <- Jet$X1 - Jet$X2`
  - o `Relation: seatpos$Seated <- seatpos$Seated / seatpos$Hi; seatpos$Arm <- seatpos$Arm / seatpos$Hi`
  - o **Variable** mit dem **höchsten VIF** aus dem Modell **entfernen** (= **«Aputation»**)
- Korrelationen treten fast zwingend auf, wenn **Vergleich zur Anzahl Beobachtungen viele erkl. Var.** vorhanden sind.
  - Setze so genannte Shrinkage-Schätzer ein wie z.B.
    - o Hauptkomponentenregression (principal component regression), ridge regression oder LASSO, elastic net.
    - o Solche Schätzer haben sich vorteilhaft bei Prognosemodellen erwiesen. Allerdings leidet Interpretierbarkeit stark.

**Strategien zur Modellbildung**

**Warum kein Modell mit allen erklärenden Variablen?**

- **Einfachheit** – gibt es mehrere Erklärungen, so sollte man das einfachste Modell wählen.
- **Reduktion der Schätzvariabilität** – unnötige erkl. Var. führen zu höheren Variabilität in Schätzung Modellkoeffizienten.
- **Bessere Interpretierbarkeit** – Multikollinearität führt zu nicht eindeutigen Lösungen.
- **Vorhersage** – weniger erkl. Var., weniger Aufwand.

**Automatisierte Variablenauswahlverfahren genügen nicht immer, weil:**

- Auswahl Variablen ist vom Zufall abhängig → neben «besten» Modell die «fast gleich guten» auch in Betracht ziehen.
- Das «beste» Modell muss nicht unbedingt alle Modellvoraussetzungen erfüllen. → **Residuenanalyse**.
- Können nicht alle nützlichen Transformationen von Beginn weg berücksichtigen wegen grosse Anzahl Möglichkeiten.
- Manchmal liefern die Verfahren Modelle, die mit dem gesicherten Fachwissen nicht übereinstimmen.
- Die «besten» Modelle können für einen Zweck sehr tauglich sein, jedoch für einen anderen unbrauchbar.

**Strategien zur Modellbildung: Grobe Skizze zum Vorgehen**

- i. **Problem verstehen** – Gibt es schon Modellansätze? Effektschätzen oder Vorhersage?
- ii. **Daten beschaffen**, kennenlernen und aufbereiten
  - o Umgang mit fehlenden Werten
  - o Bedeutung der Zahl 0 in den verschiedenen Variablen vereinheitlichen
  - o Daten gem. first-aid Transformationen behandeln, ausser es gibt triftige Gründe dagegen; Interaktionen sinnvoll?
- iii. **Erste Anpassung** (mit allen Variablen); vorzugsweise mit robusten Methoden oder mit GAM und dann Linearisieren
- iv. **Residuen-Analyse**; Tragen Daten dazu bei, das Problem zu lösen? ev. zurück nach ii. oder i.
- v. **Variablenauswahl**, ggf. Kollinearitäten behandeln
- vi. **Modelleignung klären**
  - o Residuen-Analyse mit selektierten Modell
  - o Modell/geschätzte Parameterwerte mit Fachwissen abgleichen
  - o «out-of-sample» Validierung in Betracht ziehen (Validierung mit noch nicht verwendeten Daten)

**Ursache / Kausalität:** Grundlegend für alle Wissenschaften ist die Suche nach Ursache-Wirkungs-Beziehungen.

- Aus stat. Korrelationen kann nicht auf solche Beziehungen geschlossen werden!
- Ursachen kann mit Regres.modell von Beobachtungsstudien nicht identifizieren, aber mit Daten geplanten Versuchen.
- Dennoch besteht eine wichtige Anwendung der Regression gerade darin, **Indizien** für solche Beziehungen zu sammeln.
- Die folgenden zwei Arten von Schlüssen sind üblich.
  - o Nachweis für vermutete ursächliche Wirkung der Eingangsgrösse auf die Zielgrösse, falls ein Koeffizient in einem Regressionsmodell signifikant von Null verschieden ist und eine ursächliche Wirkung der Zielgrösse auf die Eingangsvariable ausgeschlossen werden kann (Erschütterung kann Distanz zum Sprengort nicht beeinflussen!)
  - o Off Korrelation zwischen Eingangsvariable und Zielgrösse durch dritte Grösse zustande → **Scheinkorrelation**

**Additive Modelle**

- Möglichkeit, eine geeignete **nicht-lineare** Funktion  $h$  anzupassen → «Nichtlineare Regression»
- Allerdings muss eine solche Funktion vom **Fachwissen motiviert** werden. Wenn dieses Wissen fehlt, was dann?
- Möchten **funktionalen Zusammenhang**  $h$  aus Daten bestimmen, ohne spezifische Annahmen an  $h$  ausser «**Glattheit**».
- Solche Schätzverfahren für  $h$  heissen **Glätter/Nichtparametrische Schätzung** des funktionalen Zusammenhangs.

**Polynomiale Regression:** Jede stetige Fun. kann **Polynom approximieren**. Je höher Ordnung, je bessere Approximation

- `plot(accel ~ times, data = mcycle); h.new <- data.frame(times = seq(min(mcycle$times), max(mcycle$times), length = 100)); #x-Werte für Approximation MU.lm3 <- lm(accel ~ poly(times, 3), data = mcycle) #3 = Polynomgrad`
- `lines(h.new$times, predict(MU.lm3, newdata = h.new)); MU.lm6 <- lm(accel ~ poly(times, 6), data = mcycle)`
- `lines(h.new$times, predict(MU.lm6, newdata = h.new), col = 2); legend("topright", legend=c(3,6), text.col=c(1,2))`
- Polynome mit höheren Ordnungen führen zu unpassenden Schwingungen an den beiden Enden → Alternativen gefragt.
- **Zwei Konzepte**, um glatte Funktionen an Daten anzupassen → **Spline-Interpolation/lokale Regression**.

**Splines**

- Splines sind **stückweise Polynome** der Ordnung  $k$ . Die Verbindungspunkte der Stücke werden als **Knoten** bezeichnet.
- I.A. müssen Funktionswerte und die ersten  $(k - 1)$  Ableitungen an den Knoten übereinstimmen, damit der Spline eine **kontinuierliche** Funktion mit  $(k - 1)$  kontinuierlichen Ableitungen ist. Kubische Spline ( $k = 3$ ) ist meistens ausreichend.
- Um Idee in Datenanalyse zu übertragen, hilft es, Spline-Kurve als Linearkombination von  $q$  Basisfunktionen dazustellen:
- Legen  $(q + 3)$  Stützstellen fest, sodass  $x_1^* < x_2^* < \dots < x_{q+3}^*$ . Der für Berechnung relevante Bereich liegt in  $[x_2^*, x_{q+1}^*]$ .

**B-Spline-Basisfunktionen:** Kubische Spline-Kurve:  $g(x) = \beta_1 \cdot b_1^{(2)}(x) + \beta_2 \cdot b_2^{(2)}(x) + \dots + \beta_q \cdot b_q^{(2)}(x)$ .

- `library(splines); plot(accel ~ times, data = mcycle) #h.new gleich wie bei Polynom, degree = 3 → Default, kubisch`
- `MU5 <- lm(accel ~ bs(times,df = 5, degree=3), data = mcycle); lines(h.new$times, predict(MU5, newdata = h.new))`
- **#Knoten:** `h.knots <- seq(min(mcycle$times), max(mcycle$times), length=6)[-c(1,6)]; MU.k <- lm(accel ~ bs(times, knots = h.knots), data = mcycle); lines(h.new$times, predict(MU.k, newdata = h.new)) #knots = Innereknöten`

**Regressions-Splines**

- Oberer Ansatz in Regression übertragen. Dort rauschen die  $y_i$  → besser Daten zu glätten und nicht interpolieren. Deshalb
- Ein Raster von  $q$  Stützpunkten wird gewählt.
  - o Üblicherweise werden die **Stützstellen** auf **regelmässigen Gitter** über dem Wertebereich von  $x$  gewählt oder die Stellen werden durch **Quantile** definiert.
- Glattheit wird über Anzahl Stützpunkte  $q$  bestimmt. Ist  $q$  klein, ist Kurve **sehr glatt**, ist  $q$  gross, so führt zu **Interpolation**.

**Smoothing Splines** – Wie soll  $q$  (d.h. df) – und damit die Glattheit der Spline-Kurve – gewählt werden?

- Lösung: Feste, jedoch unterglättende Spline-Approximation. Bei Minimierung der Kriteriumsfunction Bestrafung für «Welligkeit» (Komplexität) hinzufügt. Anstelle von  $(y - x\beta)^T (y - x\beta)$  wird  $(y - x\beta)^T (y - x\beta) + \lambda \cdot (\text{Welligkeit})^2$  minimiert.
- Bestrafung (*Welligkeit*)<sup>2</sup> kann als quadrierte Form der Koeffizienten geschrieben werden,  $(\text{Welligkeit})^2 = \beta^T S \beta$ .
- Matrix  $S$  wird aus Basisfunktionen abgeleitet. Zielkonflikt Modellgenauigkeit/glattheit wird durch Parameter  $\lambda$  geregelt.
- Ist  $\lambda$  sehr gross → erhalten wir eine Gerade (glatteste Funktion). Ist  $\lambda = 0$  → unterglättete Smoothing-Spline-Kurve.
- Diese beiden Extreme von  $\lambda$  sind nicht praxistauglich. Es gibt ein geeignetes  $\lambda \in (0, \infty)$ , das prakt. Anforderung genügt.
- Optimale Glattheit mit optimalen Wahl  $\lambda$ , welche **Kreuzvalidierungsverfahren** basiert. ↓ ohne Angabe df ist optimal
- `library(splines); plot(accel ~ times, data = mcycle); MU.ss <- smooth.spline(x = mcycle$times, y = mcycle$accel)`
- `lines(MU.ss); MU.ss5 <- smooth.spline(x = mcycle$times, y = mcycle$accel, df = 5); lines(MU.ss5, col = 2)`

**Thin plate regression splines** (nachfolgend «**trps**» genannt) `library(mgcv); gam(y ~ s(x), data = dat)`

- Vorgestellt Ansatz hat Verbesserungspotential, da man Stützstellen wählen muss → Willkür beim Anpassungsprozess.
- Ansatz mit Basisfunktionen geht nur für **eine** erkl. Var. **Ausweg:** `trps` → **knotenfreien Basisfunktionen**, **mehr. erkl. Var.** glätten und ist «**optimal**». **Nachteil:** Schwer darstellen und nicht zu rechenintensiv Implementation haben Approx.

**Lokale Regression – LOWESS (LOcally WEighted Scatterplot Smoother)**

- LOWESS beruht auf Anwendung der einfachen Regressionsrechnung in einer lokalen Umgebung.
- Idee: kann jede Funktion  $h$  in einer kleinen Umgebung von einem vorgegebenen Punkt linear approximieren.
- Wählen ein Fenster um einen Punkt  $z_1$  aus, bei dem wir  $h(z_1)$  bestimmen wollen.
- Wählen Fenster so, dass die Funktion  $h$  möglichst gut durch eine Gerade approximierbar wird.
- Passen Gerade an Punkte an, die im Fenster liegen und machen damit Vorhersage an Stelle  $z_1$  → schätzt Wert  $h(z_1)$
- Vorgehen wird für Menge von Punkten  $z_i, i = 1, \dots, N$ , die Bereich der  $x$ -Werte möglichst gut überdecken, durchgeführt.
- Man erhält dann für die Punkte  $z_1, \dots, z_N$  die entsprechenden Vorhersagen  $\hat{h}(z_1), \dots, \hat{h}(z_N)$ .
- Um Funktion  $h$  zu visualisieren, werden Punkte durch Geradenstücke verbunden.

- Idee wird mit einer gewichteter KQ-Methode umgesetzt:  $\hat{\beta}(z_1) = \arg \min_{\beta} \sum_{i=1}^n K\left(\frac{x_i - z_1}{b}\right) (y_i - (\beta_0 + \beta_1(x_i - z_1)))^2$
- Dabei ist  $b$  **halbe Fensterbreite** (bandwidth) und  $K\left(\frac{x_i - z_1}{b}\right)$  **Gewichtsfunktion** (Kerngewicht/kernel weight) für Punkt  $i$ .
- **Offene Punkte: Wahl Fensterbreite  $b$ , Wahl Gewichtsfunktion  $K(u)$  und Behandlung Ausreissern  $y$ -Richtung.**

**Wahl der Fensterbreite  $b$ :** Wenn wir  $b$  sehr klein wählen,

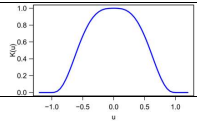
- wird der **Approximationsfehler** sehr **klein** sein – was vorteilhaft ist.
- **Nachteil grosse Varianz** der Vorhersage  $\hat{h}(z_1)$ , weil Anzahl Punkte im Fenster vergleichsweise klein sein wird.
- **Machen wir die Fensterbreite  $b$  zu gross,** dann **können wir  $\hat{h}(z_1)$  sehr genau schätzen**, d.h. mit **kleiner Varianz**,
- jedoch kann dann die **lineare Approximation ungenügend** sein, d.h. der **Approximationsfehler** wird **gross**.

Um sicherzustellen, dass immer genügend Punkte im Fenster sind, machen wir die Fensterbreite folgendermassen:

- Fensterbreite  $b$  so wählen, dass vorgegebene Anzahl  $d$  Punkte  $x_i$  im Fenster liegen. Also ist  $b$  für eine Stelle  $z_1$  der  $d$ -kleinste Wert unter  $|x_i - z_1|$ , für  $i = 1, \dots, n$ .
- Empfohlen wird für die Anzahl  $d$  der Wert  $d = \frac{2}{3} \cdot \text{Anzahl Punkte } n$ . Hilfreich ist es, einige **Werte auszuprobieren**.

**Gewichtsfunktion  $K(u)$**

- Die eigentliche Wahl der Gewichtsfunktion  $K$  ist weder theoretisch noch empirisch sehr wichtig.
- Kann  $K$  so wählen, dass z.B. die Berechnung vereinfacht wird.
- Im LOWESS/LOESS wird der «tricube kernel» benutzt:  $K(u) = (\max\{1 - |u|^3, 0\})^3$



**Umsetzung mit LOWESS / LOESS**

- $\hat{\beta}(z_1) = \arg \min_{\beta} \sum_{i=1}^n w_r(x_i) \cdot K\left(\frac{x_i - z_1}{b}\right) (y_i - (\beta_0 + \beta_1(x_i - z_1)))^2$  Robustheitsgewichte  $w_r(x_i)$  sind nur implizit bestimmt.
- Zur Erinnerung: Tukey's bisquare Gewichtsfunktion lautet  $w_r(x_i) = \left(\max\{1 - \left(\frac{r_i}{c}\right)^2, 0\}\right)^2$  mit  $r_i = \frac{y_i - \hat{h}(x_i)}{\hat{\sigma}_{MAV}}$
- $c$  kontrolliert Einfluss «schlechte» Beobachtung. Je kleiner  $c$ , desto schneller verlieren extreme Beobach. an Einfluss.
- Preis für (zu) kleine  $c$ : Effizienzverlust und grosse Standardabw. der Schätzungen → Cleveland schlägt  $c = 4.05$  vor.
- **Lokale KQ:** `scatter.smooth(P$prestige, P$income, span = 0.35, degree = 1, family = "gaussian")`
- **Robust:** `lines(loess.smooth(P$prestige, P$income, span = 0.35, degree = 1, family = "symmetric"), col = 2)`
- **Alternative:** `MU.lr <- loess(accel ~ times, data = mcycle); lines(hx$times, predict(MU.lr1g75, newdata = h.new))`
- **Default-Werte:** span = 0.75 → Fensterbreite | degree = 1 → lokal linear, Alternative = 2 → lokal quadratisch | family = "gaussian" → Gauss'sche Abweichung, Alternative "symmetric" = robuste Anpassung/langschwänzige Abweichungen

**Das additive Modell**

- Bei generalisierte additive Modell (GAM) wird (generalisierte) lineare Regressionsmodell dadurch erweitert, dass **einige** oder **alle linearen Terme** des linearen Prädiktors durch **geeignete glatte Funktionen** der Terme **ersetzt** werden.
- Ersetzung lineare Prädiktor  $\eta_i = \beta_0 + \sum_{k=1}^m x_i^{(k)} \beta_k$  durch additiven Prädiktor  $\eta_i = \sum_{k=1}^m f_k(x_i^{(k)})$
- **Einfachste**, glatte Funktionen sind **linear**. Führt zu üblichen lin. Modellen. I.A. Funktionen  $f_k(\cdot), k = 1, \dots, m$ , **unbekannt**.
- Schätzung glatte, unbekannte Funktionen führt zu Streudiagrammglätter → Smoothing-Splines oder LOWESS/LOESS
- Smoothing-Splines wird mit mehreren Variablen mit `trps` durchgeführt (R BF siehe oben → `trps`).

**LOWESS: library(gam); P.gam <- gam(prestige ~ lo(income) + lo(education), data = Prestige); summary(P.gam)**

Anova for Parametric Effects = linearer Teil

	Df	Sun Sq	Mean Sq	F value	Pr(>F)
lo(income)	1.000	14899.4	14899.4	304.20	< 2.2e-16 ***
lo(education)	1.000	8159.2	8159.2	166.58	< 2.2e-16 ***
Residuals	92.599	4535.5	49.0		

**Residuenplot 95%-Vertrauensband: par(mfrow=c(2,2)); Anova for Nonparametric Effects = nichtlinearer Teil**

`plot(P.gam, se = T); plot(P.gam, se = T, residuals = T)`

	Npar	Df	Npar F	Pr(>F)
(Intercept)				
lo(income)	4.0	5.4878	0.0005086	***
lo(education)	2.4	3.3730	0.0309103	*

**Residuen-Plot für GAM-Anpassung: par(mfrow = c(2, 2)); stats::plot.lm(P.gam) #gleich wie reg. Residuenanalyse**

- **Smoothing-Splines:** es wird kubische Splines mit Knoten an Stellen  $x_i$  eingesetzt. Steuerung **Glättungsgrad** über df.
- df = äquivalente Freiheitsgrade und ist **Mass für Glattheit** der Kurve. df entspricht der Spur der Glättungsmatrix  $S$ .
- `library(gam); P.gam <- gam(prestige ~ s(income) + s(education), data = Prestige); s(income, df = 4) #4 = Default`
- **Residuenplot 95%-Vertrauensband:** `par(mfrow=c(2,2)); plot(P.gam, se = T); plot(P.gam, se = T, residuals = T)`

**Anmerkungen GAM – Hauptnutzen GAM:** Aufdecken von **geeigneten Transformationen der erklärenden Variablen**.

- GAM schätzt diese Transformationen (d.h.  $f_k(\cdot)$ ) **nichtparametrisch** (d.h. bestimmt sie Daten gestützt).
- Analyse der notwendigen Transformationen erfolgt über partiellen Residuenplot der geschätzten Funktionen.
- `Plot.gam` zeigt für jeden additiven Term das Streudiagramm, mit Anwendung Glätter, und die resultierende glatte Kurve.

## Stand der Arbeit

SW	Vorlesung	AB/Aufgabe	Skript	ZF
1	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
2	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
3	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
4	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
5	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
6	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
7	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
8	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
9	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
10	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
11	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
12	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
13	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
14	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

Woche	Unterricht	Skript	Praktikum
Woche 1 (KW 08)	Einführung und einfache lineare Regression	Kap. 1 bis 2.4	PA 1
Woche 2 (KW 09)	Vertrauens- und Prognosebereiche; Residuen-Analyse	Kap. 2.5 bis 3.1	PA 1, PA 2
Woche 3 (KW 10)	Residuen-Analyse, Behandlung von Unzulänglichkeiten	Kap. 3.1 bis 3.4	PA 2
Woche 4 (KW 11)	Multiple lineare Regression: Die multiple Erweiterung	Kap. 4.1 und 4.2	PA 3
Woche 5 (KW 12)	Multiple lineare Regression: Schätzung und Inferenz	Kap. 4.2 und 4.5	PA 3
Woche 6 (KW 13)	Prüfen der Modelleignung	Kap. 5.1 und 5.2	PA 4
Woche 7 (KW 14)	Gewichtete Regression, einflussreiche Beobachtungen	Kap. 5.3 und 5.4	PA 4
Woche 8 (KW 15)	Robuste Anpassungsmethoden	Kap. 5.5	PA 4
Woche 9 (KW 16)	Variablenselektion, Zwischenprüfung	Kap. 6.1, 6.2/6.A	PA 5
Woche 10 (KW 17)	Variablenselektion (II)	Kap. 6.2 bis 6.3	PA 5
Woche 11 (KW 18)	Kollinearität, Strategien bei der Modellentwicklung	Kap. 6.4 und 6.5	PA 5
Woche 12 (KW 19)	Spline, lokale Regression	Kap. 7.1 bis 7.3	PA 6
Woche 13 (KW 20)	Additive Modelle, Fallstudie (nicht im Skript)	Kap. 7.4	PA 6
Woche 14 (KW 21)	Ausblick, Besprechung Zwischenprüfung, Fragestunde	Kap. 8	PA 6

### Modulaufgaben

Erlaubte Hilfsmittel für Modulendprüfung: 6 A4-Seiten Zusammenfassung auf Papier (sprich 3 A4 Seiten doppelseitig), einschliesslich R-Befehle. Es sind keine weiteren Hilfsmittel und elektronische Kommunikationsmittel wie z.B. Mobiltelefone und WLAN erlaubt.